

---

Proceedings of the  
International Congress of Mathematicians  
August 21–29, 1990, Kyoto, Japan

---





# Proceedings of the International Congress of Mathematicians



August 21-29, 1990  
Kyoto, Japan

KYOTO Volume II

The Mathematical Society of Japan



Springer-Verlag

Tokyo Berlin Heidelberg New York  
London Paris Hong Kong Barcelona  
Budapest

International Congress of Mathematicians  
August 21-29, 1990, Kyoto, Japan

*Editor:*

Ichiro Satake  
Mathematical Institute  
Faculty of Science  
Tohoku University  
Sendai 980, Japan

---

The logo for the ICM-90, designed by Kazuyoshi Aoki and Yuji Komai, symbolizes a Japanese stone lantern, the first character for Kyoto, as well as the character for  $10^{16}$ .

---

With 96 figures, including 11 halftone illustrations

ISBN 4-431-70047-1 Set (2 volumes)  
Springer-Verlag Tokyo Berlin Heidelberg New York  
ISBN 3-540-70047-1 in 2 Bänden  
Springer-Verlag Berlin Heidelberg New York Tokyo  
ISBN 0-387-70047-1 Set (2 volumes)  
Springer-Verlag New York Berlin Heidelberg Tokyo

Library of Congress Cataloging-in-Publication Data  
International Congress of Mathematicians (1990: Kyoto, Japan)  
Proceedings of the International Congress of Mathematicians, August 21-29, 1990, Kyoto /  
edited by Ichiro Satake.  
p. cm. Includes bibliographical references. ISBN 0-387-70047-1  
1. Mathematics – Congresses. I. Satake, Ichiro. 1927-. II. Title QA1.I82 1990 510-dc20  
91-4972 CIP

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in other ways, and storage in data banks.

© The Mathematical Society of Japan 1991  
Printed in Hong Kong

The use of registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typesetting: Asco Typesetting Ltd., Hong Kong; Springer Te<sub>X</sub> in-house system and typesetting output by Universitätsdruckerei H. Stürtz AG, Würzburg, Fed. Rep. of Germany  
Printing and binding: Permanent Typesetting & Printing Co., Ltd., Hong Kong

# Contents

## *Volume I*

---

Contents . . . . .	vii
Past Congresses . . . . .	ix
Past Fields Medalists and Rolf Nevanlinna Prize Winners . . . . .	xi
Organization of the Congress . . . . .	xiii
The Organizing Committee of the	
International Congress of Mathematicians 1990 . . . . .	xv
List of Donors . . . . .	xxii
Opening Ceremonies . . . . .	xxv
Closing Ceremonies . . . . .	xxxix
Scientific Program . . . . .	xliii
Invited One-Hour Addresses at the Plenary Sessions	
Invited Forty-Five Minute Addresses at the Section Meetings	
List of Participants . . . . .	li
Membership by Nationality . . . . .	lxxxvii
The Work of the Fields Medalists and the Rolf Nevanlinna Prize Winner . .	1
Invited One-Hour Addresses at the Plenary Sessions . . . . .	41
Invited Forty-Five Minute Addresses at the Section Meetings . . . . .	301
Section 1 Mathematical Logic and Foundations . . . . .	303
Section 2 Algebra . . . . .	317
Section 3 Number Theory . . . . .	403
Section 4 Geometry . . . . .	491
Section 5 Topology . . . . .	599
Section 6 Algebraic Geometry . . . . .	699
Author Index . . . . .	767

*Volume II*

---

Contents . . . . .	v
Scientific Program . . . . .	vii
Section 7 Lie Groups and Representations . . . . .	769
Section 8 Real and Complex Analysis . . . . .	847
Section 9 Operator Algebras and Functional Analysis . . . . .	969
Section 10 Probability and Mathematical Statistics . . . . .	1025
Section 11 Partial Differential Equations . . . . .	1113
Section 12 Ordinary Differential Equations and Dynamical Systems .	1235
Section 13 Mathematical Physics . . . . .	1305
Section 14 Combinatorics . . . . .	1421
Section 15 Mathematical Aspects of Computer Science . . . . .	1479
Section 16 Computational Methods . . . . .	1549
Section 17 Applications of Mathematics to the Sciences . . . . .	1607
Section 18 History, Teaching and the Nature of Mathematics . . . . .	1639
Author Index . . . . .	1683

# Scientific Program

## Invited One-Hour Addresses at the Plenary Sessions

Spencer Bloch: Algebraic $K$ -Theory, Motives, and Algebraic Cycles . . . . .	43
Stephen A. Cook: Computational Complexity of Higher Type Functions . . . . .	55
Boris L. Feigin: Conformal Field Theory and Cohomologies of the Lie Algebra of Holomorphic Vector Fields on a Complex Curve . . . . .	71
Andreas Floer: Elliptic Methods in Variational Problems . . . . .	87
Yasutaka Ihara: Braids, Galois Groups, and Some Arithmetic Functions . . . . .	99
Vaughan F. R. Jones: Von Neumann Algebras in Mathematics and Physics . . . . .	121
László Lovász: Geometric Algorithms and Algorithmic Geometry . . . . .	139
George Lusztig: Intersection Cohomology Methods in Representation Theory . . . . .	155
Andrew J. Majda: The Interaction of Nonlinear Analysis and Modern Applied Mathematics . . . . .	175
Grigorii A. Margulis: Dynamical and Ergodic Properties of Subgroup Actions on Homogeneous Spaces with Applications to Number Theory . . . . .	193
Richard B. Melrose: Pseudodifferential Operators, Corners and Singular Limits . . . . .	217
Shigefumi Mori: Birational Classification of Algebraic Threefolds . . . . .	235
Yakov G. Sinai: Hyperbolic Billiards . . . . .	249
Karen Uhlenbeck: Applications of Non-Linear Analysis in Topology . . . . .	261
Alexandre Varchenko: Multidimensional Hypergeometric Functions in Conformal Field Theory, Algebraic $K$ -Theory, Algebraic Geometry . . . . .	281

## Invited Forty-Five Minute Addresses at the Section Meetings

### Section 1: Mathematical Logic and Foundations

Ehud Hrushovski: Categorical Structures ( <i>Manuscript not submitted</i> )	
Theodore A. Slaman: Degree Structures . . . . .	303
John R. Steel: Iteration Trees ( <i>Manuscript not submitted</i> )	
Lou P. van den Dries: The Logic of Local Fields ( <i>Manuscript not submitted</i> )	
There were 15 Short Communications in this section.	

## Section 2: Algebra

Jon F. Carlson: Cohomology and Modules over Group Algebras . . . . .	317
Rostislav I. Grigorchuk: On Growth in Group Theory . . . . .	325
Craig Hunke: Absolute Integral Closure and Big Cohen-Macaulay Algebras . . . . .	339
Alexander R. Kemer: Identities of Associative Algebras . . . . .	351
Paul C. Roberts: Intersection Theory and the Homological Conjectures in Commutative Algebra . . . . .	361
Klaus W. Roggenkamp: The Isomorphism Problem for Integral Group Rings of Finite Groups . . . . .	369
Robert W. Thomason: The Local to Global Principle in Algebraic $K$ -Theory . . . . .	381
Efim I. Zelmanov: On the Restricted Burnside Problem . . . . .	395

There were 59 Short Communications in this section.

## Section 3: Number Theory

Henri Gillet: A Riemann-Roch Theorem in Arithmetic Geometry . . . . .	403
Martin N. Huxley: Area, Lattice Points and Exponential Sums . . . . .	413
Kazuya Kato: Generalized Class Field Theory . . . . .	419
Victor Alecsandrovich Kolyvagin: On the Mordell-Weil Group and the Shafarevich-Tate Group of Modular Elliptic Curves . . . . .	429
Gérard Laumon: La Transformation de Fourier Géométrique et ses Applications . . . . .	437
Yuri Nesterenko: Algebraic Independence of Values of Analytic Functions	447
Peter C. Sarnak: Diophantine Problems and Linear Groups . . . . .	459
Tetsuji Shioda: Theory of Mordell-Weil Lattices . . . . .	473

There were 29 Short Communications in this section.

## Section 4: Geometry

Kenji Fukaya: Collapsing Riemannian Manifolds and Its Applications . . . . .	491
Etienne Ghys: Le Cercle à l'Infini des Surfaces à Courbure Négative . . . . .	501
Karsten Grove: Metric and Topological Measurements of Manifolds . . . . .	511
Helmut Hofer: Symplectic Invariants . . . . .	521
Peter B. Kronheimer: Embedded Surfaces in 4-Manifolds . . . . .	529
Dusa McDuff: Symplectic 4-Manifolds . . . . .	541
John J. Millson: Rational Homotopy Theory and Deformation Problems from Algebraic Geometry . . . . .	549
Eugenii I. Shustin: Geometry of Discriminant and Topology of Algebraic Curves . . . . .	559
Joseph H. M. Steenbrink: Applications of Hodge Theory to Singularities .	569
Toshikazu Sunada: (with M. Nishio) Trace Formulae in Spectral Geometry	577
Gang Tian: Kähler-Einstein Metrics on Algebraic Manifolds . . . . .	587

There were 38 Short Communications in this section.

## Section 5: Topology

Marcel A. Bökstedt: Algebraic $K$ -Theory of Spaces and the Novikov Conjecture ( <i>Manuscript not submitted</i> )	
Francis Bonahon: Ensembles Limites et Applications . . . . .	599
David Gabai: Foliations and 3-Manifolds . . . . .	609
Thomas G. Goodwillie: The Differential Calculus of Homotopy Functors . . . . .	621
Cameron McA. Gordon: Dehn Surgery on Knots . . . . .	631
Kiyoshi Igusa: Parametrized Morse Theory and Its Applications . . . . .	643
Lowell E. Jones: (with F. T. Farrell) Rigidity in Geometry and Topology . . . . .	653
Shigeyuki Morita: Mapping Class Groups of Surfaces and Three-Dimensional Manifolds . . . . .	665
Henri Moscovici: Cyclic Cohomology and Invariants of Multiply Connected Manifolds . . . . .	675
Vladimir G. Turaev: State Sum Models in Low-Dimensional Topology . . . . .	689

There were 43 Short Communications in this section.

## Section 6: Algebraic Geometry

Yujirō Kawamata: Canonical and Minimal Models of Algebraic Varieties . . . . .	699
János Kollar: Flip and Flop . . . . .	709
Robert K. Lazarsfeld: Linear Series on Algebraic Varieties . . . . .	715
Morihiko Saito: Mixed Hodge Modules and Applications . . . . .	725
Leslie Saper: $L_2$ -Cohomology of Algebraic Varieties . . . . .	735
Carlos T. Simpson: Nonabelian Hodge Theory . . . . .	747
Paul Vojta: Arithmetic and Hyperbolic Geometry . . . . .	757

There were 25 Short Communications in this section.

## Section 7: Lie Groups and Representations

Dan Barbasch: Unipotent Representations for Real Reductive Groups . . . . .	769
Günter Harder: Eisenstein Cohomology of Arithmetic Groups and Its Applications to Number Theory . . . . .	779
Masaki Kashiwara: Crystallizing the $q$ -Analogue of Universal Enveloping Algebras . . . . .	791
Olivier Mathieu: Classification of Simple Graded Lie Algebras of Finite Growth . . . . .	799
Toshihiko Matsuki: Orbits on Flag Manifolds . . . . .	807
Colette Moeglin: Sur les Formes Automorphes de Carré Intégrable . . . . .	815
Gopal Prasad: Semi-simple Groups and Arithmetic Subgroups . . . . .	821
Stephen Rallis: Poles of Standard $L$ Functions . . . . .	833

There were 23 Short Communications in this section.

## Section 8: Real and Complex Analysis

Eric Bedford: Iteration of Polynomial Automorphisms of $C^2$ . . . . .	847
Michael Christ: Precise Analysis of $\bar{\partial}_b$ and $\bar{\partial}$ on Domains of Finite Type in $C^2$ . . . . .	859
Ronald R. Coifman: Adapted Multiresolution Analysis, Computation, Signal Processing and Operator Theory . . . . .	879
Curt McMullen: Rational Maps and Kleinian Groups . . . . .	889
Takafumi Murai: Analytic Capacity for Arcs . . . . .	901
Takeo Ohsawa: Recent Applications of $L^2$ Estimates for the Operator $\bar{\partial}$ .	913
David Preiss: Differentiability and Measures in Banach Spaces . . . . .	923
Kyoji Saito: The Limit Element in the Configuration Algebra for a Discrete Group: A précis . . . . .	931
Nessim Sibony: Some Recent Results on Weakly Pseudoconvex Domains	943
Nicholas Th. Varopoulos: Analysis and Geometry on Groups . . . . .	951
Alexander L. Volberg: Asymptotically Holomorphic Functions and Certain of Their Applications . . . . .	959

There were 73 Short Communications in this section.

## Section 9: Operator Algebras and Functional Analysis

Joachim Cuntz: Cyclic Cohomology and $K$ -Homology . . . . .	969
Adrian Ocneanu: Quantum Symmetry and Classification of Subfactors <i>(Manuscript not submitted)</i>	
Michael V. Pimsner: $K$ -Theory for Groups Acting on Trees . . . . .	979
Sorin Teodor Popa: Subfactors and Classification in von Neumann Algebras . . . . .	987
Georges Skandalis: Operator Algebras and Duality . . . . .	997
Michel Talagrand: Some Isoperimetric Inequalities and Their Applications . . . . .	1011

There were 52 Short Communications in this section.

## Section 10: Probability and Mathematical Statistics

Martin T. Barlow: Random Walks and Diffusions on Fractals . . . . .	1025
Persi Diaconis: Applications of Group Representations to Statistical Problems . . . . .	1037
Roland L. Dobrushin: Large Deviation of Gibbsian Fields <i>(Manuscript not submitted)</i>	
Richard Durrett: Stochastic Models of Growth and Competition . . . . .	1049
Hillel Furstenberg: Recurrent Ergodic Structures and Ramsey Theory .	1057
Shinichi Kotani: Random Schrödinger Operators . . . . .	1071
Shigeo Kusuoka: De Rham Cohomology of Wiener-Riemannian Manifolds . . . . .	1075
Lucien M. Le Cam ( <i>not delivered at the Congress</i> ): Some Recent Results in the Asymptotic Theory of Statistical Estimation . . . . .	1083

Stanislav A. Molchanov: Localization and Intermittency: New Results . . . . .	1091
Marc Yor: The Laws of Some Brownian Functionals . . . . .	1105

There were 37 Short Communications in this section.

### **Section 11: Partial Differential Equations**

Demetrios Christodoulou: The Stability of Minkowski Spacetime . . . . .	1113
Jean-Michel Coron: Harmonic Maps with Values into Spheres . . . . .	1123
Matthias Günther: Isometric Embeddings of Riemannian Manifolds . . . . .	1137
Mitsuru Ikawa: On Scattering by Obstacles . . . . .	1145
Gilles Lebeau: Interaction des Singularités Faibles pour les Équations d'Ondes Semi-linéaires . . . . .	1155
Fang Hua Lin: Static and Moving Defects in Liquid Crystals . . . . .	1165
Pierre-Louis Lions: On Kinetic Equations . . . . .	1173
Pierre Schapira: Sheaf Theory for Partial Differential Equations . . . . .	1187
Michael Struwe: The Evolution of Harmonic Maps . . . . .	1197
Kanehisa Takasaki: Integrable Systems in Gauge Theory, Kähler Geometry and Super KP Hierarchy – Symmetries and Algebraic Point of View . . . . .	1205
Luc Tartar: H-Measures and Applications . . . . .	1215
Michael E. Taylor: Microlocal Analysis in Spectral and Scattering Theory and Index Theory . . . . .	1225

There were 48 Short Communications in this section.

### **Section 12: Ordinary Differential Equations and Dynamical Systems**

César Camacho: Problems on Limit Sets of Foliations on Complex Projective Spaces . . . . .	1235
Lennart Carleson: The Dynamics of Non-uniformly Hyperbolic Systems in Two Variables . . . . .	1241
Jean P. Ecalle: The Acceleration Operators and Their Applications to Differential Equations, Quasianalytic Functions, and the Constructive Proof of Dulac's Conjecture . . . . .	1249
Ju. S. Il'yashenko: Finiteness Theorems for Limit Cycles . . . . .	1259
Anatoly I. Neishtadt: Averaging and Passage Through Resonances . . . . .	1271
Sheldon E. Newhouse: Entropy in Smooth Dynamical Systems . . . . .	1285
Mary Rees: Combinatorial Models Illustrating Variation of Dynamics in Families of Rational Maps . . . . .	1295
Jean-Christophe Yoccoz: Optimal Arithmetical Conditions in Some Small Divisors Theorems ( <i>Manuscript not submitted</i> )	

There were 41 Short Communications in this section.

### Section 13: Mathematical Physics

R. J. Baxter: Hyperelliptic Function Parametrization for the Chiral Potts Model . . . . .	1305
Sergio Doplicher: Abstract Compact Group Duals, Operator Algebras and Quantum Field Theory . . . . .	1319
Joel Feldman: Introduction to Constructive Quantum Field Theory . . . . .	1335
Michio Jimbo: Solvable Lattice Models and Quantum Groups . . . . .	1343
Igor Krichever: The Periodic Problems for Two-Dimensional Integrable Systems . . . . .	1353
Antti Kupiainen: Renormalization Group and Random Systems . . . . .	1363
Nicolai Reshetikhin: Invariants of Links and 3-Manifolds Related to Quantum Groups . . . . .	1373
Albert Schwarz: Geometry of Fermionic String . . . . .	1377
Graeme Segal: Geometric Aspects of Quantum Field Theory . . . . .	1387
I. M. Sigal: Quantum Mechanics of Many-Particle Systems . . . . .	1397
Akihiro Tsuchiya: Moduli of Stable Curves, Conformal Field Theory and Affine Lie Algebras . . . . .	1409
Stanisław L. Woronowicz: Noncompact Quantum Groups <i>(Manuscript not submitted)</i>	

There were 39 Short Communications in this section.

### Section 14: Combinatorics

Noga Alon: Non-Constructive Proofs in Combinatorics . . . . .	1421
Peter J. Cameron: Infinite Permutation Groups in Enumeration and Model Theory . . . . .	1431
Alexander A. Ivanov: Geometric Presentations of Groups with an Application to the Monster . . . . .	1443
Vojtech Rödl: Some Developments in Ramsey Theory . . . . .	1455
Eva Tardos: Strongly Polynomial and Combinatorial Algorithms in Optimization . . . . .	1467
Carsten Thomassen: Graphs, Random Walks and Electrical Networks <i>(Manuscript not submitted)</i>	

There were 25 Short Communications in this section.

### Section 15: Mathematical Aspects of Computer Science

László Babai: Computational Complexity in Finite Groups . . . . .	1479
Lenore Blum: A Theory of Computation and Complexity over the Real Numbers . . . . .	1491
Alexandre L. Chistov: Efficient Factoring Polynomials over Local Fields and Its Applications . . . . .	1509
Shafi Goldwasser: Interactive Proofs and Applications . . . . .	1521
Avi Wigderson: Information Theoretic Reasons for Computational Difficulty . . . . .	1537

There were 5 Short Communications in this section.

**Section 16: Computational Methods**

Ami Harten: Recent Developments in Shock-Capturing Schemes . . . . .	1549
William M. Kahan: Paradoxes in Our Concepts of Accuracy <i>(Manuscript not submitted)</i>	
Alexander V. Karzanov: Undirected Multiflow Problems and Related Topics – Some Recent Developments and Results . . . . .	1561
Robert Krasny: Computing Vortex Sheet Motion . . . . .	1573
Masatake Mori: Developments in the Double Exponential Formulas for Numerical Integration . . . . .	1585
James Renegar: Computational Complexity of Solving Real Algebraic Formulae . . . . .	1595

There were 26 Short Communications in this section.

**Section 17: Applications of Mathematics to the Sciences**

Philip Holmes: ( <i>with G. Berkooz and J. L. Lumley</i> ) Turbulence, Dynamical Systems and the Unreasonable Effectiveness of Empirical Eigenfunctions . . . . .	1607
Yves F. Meyer: Wavelets and Applications . . . . .	1619
Masayasu Mimura: Pattern Formation in Reaction-Diffusion Systems . .	1627

There were 18 Short Communications in this section.

**Section 18: History, Teaching and the Nature of Mathematics**

Annick M. Horiuchi: The Development of Algebraic Methods of Problem-Solving in Japan in the Late Seventeenth and the Early Eighteenth Centuries . . . . .	1639
Jesper Lützen: The Birth of Spectral Theory – Joseph Liouville’s Contributions . . . . .	1651
Yuri Ivanovich Manin ( <i>delivered by Barry Mazur</i> ): Mathematics as Metaphor . . . . .	1665
Haruo Murakami: Teaching Mathematics to Students Not Majoring in Mathematics – Present Situation and Future Prospects – . . . . .	1673

There were 22 Short Communications in this section.

In addition there were 2 Short Communications in the Post Deadline Session.



# Unipotent Representations for Real Reductive Groups

Dan Barbasch

Cornell University, Ithaca, NY 14853 and  
Rutgers University, New Brunswick, NJ 08903, USA

## 1. Introduction

The classification of the unitary dual of a real reductive group is an important problem in Representation theory. Typically one proceeds as follows. Given the Levi component of a parabolic subgroup, there are certain constructions such as unitary induction, derived functors and complementary series that preserve unitarity. Thus, given a group  $G$  it is reasonable to ask for a set  $\mathcal{U}(G)$  such that the unitary dual is obtained by *unitarity preserving constructions* from all such sets on Levi components of the group. Such a set is provided in the case of integral infinitesimal character by the *special unipotent representations*. In particular the question of the unitarity of this set of representations arises.

The purpose of this lecture is to describe the proof of the unitarity of a particular subset of special unipotent representations. The precise result is as follows.

**Theorem 1.1.** *The representations in the Arthur L-packet in the special unipotent case (see Section 2) for the classical groups are unitary.*

The main idea of the proof of this theorem is the same as for the corresponding fact for the complex classical groups in Section 10 of [B]. The question of unitarity of complicated representations on a group is reduced to the same question for simpler representations on larger groups. The techniques are general, so that they apply to all special unipotent representations. For (serious) technical reasons, I can only carry out the proof for the aforementioned case.

Many of the more abstract results about unipotent representations can be phrased in a more general setting. This work, joint with J. Adams and D. Vogan will appear in [ABV].

In the case of a split classical group, [M] and [MW] have shown that the spherical unipotent representations actually occur in the residual spectrum.

In Section 2, I introduce the packet  $\Pi({}^L\mathcal{O})$  of special unipotent representations and give a definition of the Arthur L-packet. Section 3 is devoted to a computation of the size of  $\Pi({}^L\mathcal{O})$ . Combined with the constructions in Section 4, this gives all the special unipotent representations. Theorem 4.4 states that for the unitarity

of the representations in the Arthur L-packet it is enough to consider only *smoothly cuspidal* (definition at the start of Section 4) nilpotent orbits. The main result is contained in Theorem 5.3 where  $\Pi(^L\mathcal{O})$  is broken up according to WF-sets in the Lie algebras of the group and the dual group. In Section 6, I sketch the proof of the unitarity of the spherical representations in the Arthur L-packet for  $Sp(2n, \mathbb{R})$ . The idea is the following. We do an induction on how close  $\pi$  is to being obtained by *unitarity preserving functors* (unitary induction from a real parabolic subalgebra and derived functor construction in the appropriate range) from a strictly smaller parabolic subalgebra. Let  $g(n)$  be the Lie algebra  $sp(2n, \mathbb{R})$  of rank  $n$  and  $\pi(^L\mathcal{O})$  be the spherical representation in  $\Pi(^L\mathcal{O})$ . Assume that it cannot be obtained by *unitarity preserving functors* from a strictly smaller parabolic subalgebra. Then, for a well chosen value  $r$ , we embed  $g(n) \times gl(r+1)$  as a Levi component of a real parabolic subalgebra  $p(\mathbb{R})$  in  $g(n+r+1)$ . We then form  $I(\pi) = \text{Ind}_{P(\mathbb{R})}^{G(\mathbb{R})}[\Pi(^L\mathcal{O}) \otimes \text{Triv}]$ . Since  $\Pi(^L\mathcal{O})$  is hermitian, the induced representation inherits an *induced hermitian form*. On the other hand, it decomposes into a sum of two irreducible representations which can be seen to arise *earlier* than  $\pi(\mathcal{O})$  in the induction. To finish the proof we need to see that the two factors have the same signature in the *induced form* of  $I(\pi)$ . This is a simple direct calculation of a signature on one K-type in  $I(\pi)$  and  $\Pi(^L\mathcal{O})$ .

Unless mentioned otherwise, I will use the following notation. The complex group  $Sp(2n, \mathbb{C})$  and its Lie algebra will be denoted by  $G$  and  $g$ . The real form  $Sp(2n, \mathbb{R})$  will be  $G(\mathbb{R})$  and its real Lie algebra  $g(\mathbb{R})$ . The dual Lie algebra  $so(2n+1)$  will be denoted  ${}^Lg$ , its various real forms  $so(p, q)$  by  $\check{g}(\mathbb{R})$ , and the corresponding groups by  ${}^LG^0$  and  $\check{G}(\mathbb{R})$ .

Complex nilpotent orbits in  $g$  will be denoted by  $\mathcal{O}$ , in the dual  ${}^Lg$  by  ${}^L\mathcal{O}$ .

$\theta$  will denote the complexification of the Cartan involution for  $Sp(2n, \mathbb{R})$ . Similarly  $\check{\theta}$  will be the complexified Cartan involution of one of the real forms of  ${}^Lg$ . The corresponding decompositions are  $g = \mathfrak{k} + \mathfrak{s}$  and  $\check{g} = \check{\mathfrak{k}} + \check{\mathfrak{s}}$ .

Real forms of the orbit  $\mathcal{O}$  will be denoted by  $\mathcal{O}(\mathbb{R})$ . Similarly for real forms of  ${}^L\mathcal{O}$ . By [S], (see also [V4] Section 4), real forms of  $\mathcal{O}$  are in 1–1 correspondence with nilpotent orbits of  $K_c$  on  $\mathfrak{s}$ . We will use this identification without explicit mention.

*Acknowledgements.* This research was partially supported by NSF grant DMS–8803500. I also wish to thank David Vogan for several very valuable comments on the material of this talk.

## 2. Special Unipotent Representations

In [A1] and [A2], Arthur introduces a Langlands parameter attached to a homomorphism

$$\psi : W_{\mathbb{R}} \times SL(2, \mathbb{C}) \longrightarrow {}^LG, \quad (2.1)$$

where  $W_{\mathbb{R}}$  is the Weil group for  $\mathbb{R}$ . He conjectures that there should be an entire packet  $\Pi(\psi)$  of representations attached to  $\psi$ , satisfying various properties relevant to the study of automorphic forms via the trace formula.

Because the group we are dealing with is split,  ${}^L G$  is a direct product. Therefore maps as in (2.1) are determined by homomorphisms into  ${}^L G^0$ , the connected group. We will only consider the case of  $\psi$ 's that are trivial on the identity component of  $W_{\mathbb{R}}$  and such that  $\psi|_{SL(2,\mathbb{C})}$  determines an even nilpotent orbit in the dual algebra. This case is called *special unipotent*.

Let  ${}^L \mathcal{O}$  be the even nilpotent orbit determined by  $\psi$  and  $\mathcal{O}$  be the dual orbit in  $g$  in the sense of [L]. Let  $E, H, F$ , be the standard generators of  $sl(2)$  and write

$${}^L e = d\psi(E), \quad {}^L h = d\psi(H), \quad {}^L f = d\psi(F). \quad (2.2)$$

The infinitesimal character of the Arthur parameter is given by

$$\lambda(\psi) = \lambda = \frac{1}{2} d\psi(H). \quad (2.3)$$

This depends only on  $\psi|_{SL(2,\mathbb{C})}$ .

**Definition 2.4.** The *extended packet*  $\Pi({}^L \mathcal{O})$  corresponding to  $\psi$  is the set of irreducible representations with annihilator in the universal enveloping algebra equal to the *special unipotent primitive ideal* attached to  ${}^L \mathcal{O}$ .

In particular,  $\overline{\text{Ad } G \cdot WF(\pi)} = \overline{\mathcal{O}}$ , for  $\pi \in \Pi({}^L \mathcal{O})$  as well as  $\overline{\text{Ad } G \cdot \mathcal{A}(\pi)} = \overline{\mathcal{O}}$ , where  $\mathcal{A}$  denotes the associated variety defined in [V4].  $\Pi({}^L \mathcal{O})$  corresponds to the union of the  $\Pi(\psi)$  giving the same nilpotent orbit  ${}^L \mathcal{O}$ .

If  $j \in W_{\mathbb{R}}$  is the element representing complex conjugation, then let  $m = \psi(j)$ . This is an element of order 2 in  $C_G({}^L e, {}^L h, {}^L f)$ . Then

$$\check{\theta}_m = \text{Ad}(m \cdot e^{i\pi\lambda}) \quad (2.5)$$

is an involution. Let  $\check{G}_m(\mathbb{R})$  be the real form of  ${}^L G$  with  $\check{\theta} = \check{\theta}_m$ . The parabolic subalgebra determined by  $\lambda$  is  $\check{\theta}$ -stable. Denote it by  $\check{\mathfrak{p}}(\lambda)$  and fix a  $\check{\theta}$ -stable Cartan subalgebra  $\check{\mathfrak{h}}$  containing  $\lambda$  and a Borel subalgebra  $\check{\mathfrak{b}} \subset \check{\mathfrak{h}}$  for which  $\lambda$  is dominant. Then  ${}^L \mathcal{O}$  meets  $\check{\mathfrak{s}} \cap \check{\mathfrak{n}}(\lambda)$ , so this also determines a real form  ${}^L \mathcal{O}(\mathbb{R})$  of  ${}^L \mathcal{O}$ .

The representations in the Arthur L-packet attached to  $\psi$  can then be constructed as follows.

Let  $\check{\mathcal{L}}(\psi) = \mathcal{A}_{\check{\mathfrak{p}}}$  be the irreducible representations (the Levi component may be disconnected) with infinitesimal character  $\check{\varrho}$  which are obtained by derived functor construction from a character of  $\check{\mathfrak{p}}(\lambda)$ , as in [V1, VZ]. [V3] attaches to each such  $\check{\mathcal{L}}(\psi)$  a standard modules  $X_{\text{reg}}(\psi)$  with regular integral infinitesimal character on a quasisplit real form  $G_{\text{qs}}(\mathbb{R})$  of  $G$ . Let  $X(\psi)$  be the standard module obtained from  $X_{\text{reg}}(\psi)$  by applying the translation functor to infinitesimal character  $\lambda$ . Then the Arthur L-packet is the set of irreducible quotients of the  $X(\psi)$ . This coincides with the usual definition of L-packet used in [A1] and [A2].

*Example.* Let  ${}^L \mathcal{O}$  be the nilpotent orbit with Jordan blocks  $31^2$  in  $so(5)$ . In the standard coordinates coming from embedding  $sp(2n)$  in  $gl(2n)$  the infinitesimal character is  $(1, 0)$  and  $\mathcal{O}$  is the nilpotent with Jordan blocks  $2^2$ . The centralizer of

the Lie triple corresponding to  ${}^L\mathcal{O}$  is isomorphic to  $S[O(1) \times O(2)]$ , so it has three conjugacy classes of elements of order 2. They correspond to two real orbits in  $so(3, 2)$  and one in  $so(4, 1)$ . The most split Cartan subgroup of  $Sp(4)$  is of the form  $MA$  where  $M \cong \mathbb{Z}_2 \times \mathbb{Z}_2$ . If we assume that the first  $\mathbb{Z}_2$  corresponds to the coordinate containing the 1 in the infinitesimal character, then the three principal series have  $M$ -characters  $\text{Triv} \otimes \text{Triv}$ ,  $\text{Sgn} \otimes \text{Sgn}$  and  $\text{Triv} \otimes \text{Sgn}$ . The Arthur L-packet of the first one contains one irreducible representation namely the spherical, the other ones two irreducible representations each.

### 3. The Coherent Continuation Representation

The results in this section hold for general reductive groups.

Fix a regular integral infinitesimal character  $\chi_{\text{reg}}$ . Denote by  $\mathcal{G}(\chi_{\text{reg}})$  the Grothendieck group of the category of  $(\mathfrak{g}, K)$  modules with infinitesimal character  $\chi_{\text{reg}}$ . Recall from [V1] that there is an action of the Weyl group on  $\mathcal{G}(\chi_{\text{reg}})$ , called the *coherent continuation action*. Then  $\mathcal{G}(\chi_{\text{reg}})$  decomposes into a direct sum according to blocks  $\mathcal{B}$ ,

$$\mathcal{G}(\chi_{\text{reg}}) = \bigoplus \mathcal{G}_{\mathcal{B}}(\chi_{\text{reg}}). \quad (3.1)$$

Let  $\mathfrak{h}_a \subset \mathfrak{g}$  be an abstract Cartan subalgebra and let  $\Pi_a$  be a set of (abstract) simple roots. For each irreducible representation  $\mathcal{L}(\gamma)$ , denote by  $\tau(\gamma)$  the tau-invariant as defined in [V1]. Given a block  $\mathcal{B}$  and disjoint orthogonal sets  $S_1, S_2 \subset \Pi_a$ , define

$$\mathcal{B}(S_1, S_2) = \{\gamma \in \mathcal{B} | S_1 \subset \tau(\gamma), S_2 \cap \tau(\gamma) = \emptyset\}. \quad (3.2)$$

If in addition we are given a nilpotent orbit  $\mathcal{O} \subset \mathfrak{g}$ , we can also define

$$\mathcal{B}(S_1, S_2, \mathcal{O}) = \{\gamma \in \mathcal{B}(S_1, S_2) | WF(\mathcal{L}(\gamma)) \subset \overline{\mathcal{O}}\}. \quad (3.3)$$

Consider the case of  $\mathfrak{g}$  viewed as a real Lie algebra. Then the case  $S_1, S_2 = \emptyset$  is called the double cone  $\mathcal{C}(\mathcal{O})$ . The double cell corresponding to  $\mathcal{O}$  will be denoted  $\overline{\mathcal{C}}(\mathcal{O})$ .

Let  $W_i = W(S_i)$ , and define

$$\begin{aligned} m_S(\sigma) &= [\sigma : \text{Ind}_{W_1 \times W_2}^W (\text{Sgn} \otimes \text{Triv})], \\ m_{\mathcal{B}}(\sigma) &= [\sigma : \mathcal{G}_{\mathcal{B}}(\chi_{\text{reg}})]. \end{aligned} \quad (3.4)$$

**Theorem 3.5** (Vogan).

$$|\mathcal{B}(S_1, S_2, \mathcal{O})| = \sum_{\sigma \otimes \sigma \in \mathcal{C}(\mathcal{O})} m_{\mathcal{B}}(\sigma) m_S(\sigma).$$

A sketch of the proof can be found in [BSS].

Recall  $\lambda = \lambda(\psi)$  from Section 1. Then  $\lambda$  defines a set  $S_2$  by

$$S_2 = S(\psi) = \{\alpha \in \Pi_a | (\alpha, \lambda) = 0\}. \quad (3.6)$$

Then

$$\Pi({}^L\mathcal{O}) = \bigcup_{\mathcal{B}} \mathcal{B}(\emptyset, S(\psi), \mathcal{O}). \quad (3.7)$$

In the classical groups case,  $m_{\mathcal{B}}(\sigma)$  is explicitly computable. For the special unipotent case,  $m_S(\sigma)$  equals 0 except for the representations occurring in the corresponding left cell  $\overline{\mathcal{C}}^L(\mathcal{O})$  when it is 1 (see [BV]).

**Theorem 3.8.**

$$|\Pi({}^L\mathcal{O})| = \sum_{\mathcal{B}} \sum_{\sigma \otimes \sigma \in \overline{\mathcal{C}}^L(\mathcal{O})} m_{\mathcal{B}}(\sigma).$$

## 4. Induction from Parabolic Subgroups

Recall that a parabolic subalgebra  $\mathfrak{p} = \mathfrak{m} + \mathfrak{n}$  is called *real for*  $\mathfrak{g}(\mathbb{R})$  if  $\theta(\mathfrak{p}) = \bar{\mathfrak{p}}$ ,  *$\theta$ -stable for*  $\mathfrak{g}(\mathbb{R})$ , if  $\theta(\mathfrak{p}) = \mathfrak{p}$ .

A nilpotent  $\mathcal{O}$  is called *induced from*  $\mathcal{O}_{\mathfrak{m}}$  if  $\mathcal{O} \cap (\mathcal{O}_{\mathfrak{m}} + \mathfrak{n})$  is dense in  $\mathcal{O}_{\mathfrak{m}} + \mathfrak{n}$ . It is called *smoothly induced* if it satisfies Hypotheses 8A and 8B in [BV].

Let  $\mathfrak{p}(\xi) \subset \mathfrak{g}$  be real or  $\theta$ -stable. Let  $\mathcal{O}_{\mathfrak{m}}(\mathbb{R}) \subset \mathfrak{m}(\mathbb{R})$  be a real form of  $\mathcal{O}_{\mathfrak{m}}$ . We define the *real induced set* from  $\mathcal{O}_{\mathfrak{m}}(\mathbb{R})$  to be

$$\text{ind}_{\mathfrak{p}}^{\mathfrak{g}}[\mathcal{O}_{\mathfrak{m}}(\mathbb{R})] = \overline{\bigcup_{t>0} \text{Ad } G(\mathbb{R}) \cdot (t\xi + \mathcal{O}_{\mathfrak{m}}(\mathbb{R}))} \setminus \bigcup_{t>0} \text{Ad } G(\mathbb{R}) \cdot (t\xi + \mathcal{O}_{\mathfrak{m}}(\mathbb{R})). \quad (4.1)$$

In the case of  $\mathfrak{p}$  real, this set is the closure of  $\text{Ad } G[\mathcal{O}_{\mathfrak{m}}(\mathbb{R}) + \mathfrak{n}(\mathbb{R})]$ .

In the case of  $\theta$ -stable parabolic subalgebra in a classical group, this can be computed using unpublished results of D. Peterson. It is the closure of a single real orbit.

We give two constructions that yield special unipotent representations from special unipotent representations on Levi subgroups of proper parabolic subalgebras.

**Induction.** Let  $\mathfrak{p}(\xi)$  be real or  $\theta$ -stable. Assume that  ${}^L\mathcal{O}_{\mathfrak{m}}$  and  ${}^L\mathcal{O}$  are such that

$$\lambda|_{[\mathfrak{m}, \mathfrak{m}]} = \lambda_{\mathfrak{m}}, \quad \mathcal{O} = \text{ind}_{\mathfrak{p}}^{\mathfrak{g}}[\mathcal{O}_{\mathfrak{m}}]. \quad (4.2)$$

Denote by  $\lambda_c$  the character of  $\mathfrak{m}$  determined by  $\lambda$  on the center.

**I. Let  $\mathfrak{p}(\xi)$  be real.** Then

$$\text{Ind}_{P(\mathbb{R})}^{G(\mathbb{R})}[\pi_{\mathfrak{m}} \otimes \lambda_c]$$

has composition series formed of irreducible representations in  $\Pi({}^L\mathcal{O})$  only.

**II. Let  $\mathfrak{p}(\xi)$  be  $\theta$ -stable.** Assume that  $\pi \in \Pi(\mathcal{O}_{\mathfrak{m}})$  is such that  $WF(\pi_{\mathfrak{m}})$  contains an orbit  $\mathcal{O}_{\mathfrak{m}}(\mathbb{R})$  such that  $\text{ind}_{\mathfrak{p}(\mathbb{R})}^{\mathfrak{g}(\mathbb{R})}[\mathcal{O}_{\mathfrak{m}}]$  meets the orbit  $\mathcal{O}$ . Then

$$\mathcal{R}_\mathfrak{p}^i[\pi_\mathfrak{m} \otimes \lambda_c]$$

are not all zero and their composition series consists of irreducible representations in  $\Pi({}^L\mathcal{O})$ .

**Coinduction.** Let  $\check{\mathfrak{p}} = \check{\mathfrak{m}} + \check{\mathfrak{n}}$  be a parabolic subalgebra,  $\check{\mathfrak{o}}$ -stable for some real form of  $\check{\mathfrak{g}}$ , such that  $\check{\mathfrak{p}}(\lambda) \subset \check{\mathfrak{p}}$ . Let  $\mathfrak{p}$  be the dual real parabolic subalgebra. Let  $\pi_\mathfrak{m} \in \Pi(\mathcal{O}_\mathfrak{m})$  be such that  $\text{ind}_{\check{\mathfrak{p}}}^{\check{\mathfrak{g}}} [{}^L\mathcal{O}_\mathfrak{m}(\mathbb{R})]$  meets the orbit  ${}^L\mathcal{O}$ . Then the irreducible representation

$$\check{\mathcal{R}}_\mathfrak{p} = P_\lambda([\mathcal{R}_{\check{\mathfrak{p}}}(\check{\pi}_\mathfrak{m})])^\vee \quad (4.3)$$

is in  $\Pi({}^L\mathcal{O})$ . Recall that  $\mathcal{O}$  is dual in the sense of [L] to the even nilpotent orbit  ${}^L\mathcal{O}$ . The same holds for  $\mathcal{O}_\mathfrak{m}$  and  ${}^L\mathcal{O}_\mathfrak{m}$ .

**Theorem 4.4.** *Assume that  $\mathcal{O}$  is smoothly induced from  $\mathcal{O}_\mathfrak{m}$ . Let  $\pi$  be in the Arthur L-packet determined by  ${}^L\mathcal{O}$ . Then there is  $\pi_\mathfrak{m}$  in the Arthur L-packet determined by  ${}^L\mathcal{O}_\mathfrak{m}$  such that  $\pi$  is the lowest K-type subquotient of  $\text{Ind}_{P(\mathbb{R})}^{G(\mathbb{R})}[\pi_\mathfrak{m}]$ .*

In view of this, we say a nilpotent orbit is *smoothly cuspidal* if it is not smoothly induced from any proper parabolic subalgebra. They are as follows.

*Type B.* The largest Jordan block is odd size and occurs an odd number of times. All smaller odd sizes occur, and they each occur an even number of times.

*Type C.* The largest Jordan block is even size and occurs an even number of times. All smaller even sizes occur and they each occur an even number of times.

*Type D.* The largest Jordan block is odd size and occurs an even number of times. All smaller odd sizes occur, and they each occur an even number of times.

For the unitarity of representations in the Arthur L-packet, we can restrict to these cases.

## 5. The Main Result

The main result concerns the  $WF$ -set of the special unipotent representations. It is used to compute the decomposition of unitarily induced representations.

Recall the quotient of the component group  $\overline{A}(\mathcal{O})$  from [L]. For a real form  $\mathcal{O}(\mathbb{R})$ , let  $A(\mathcal{O}(\mathbb{R}))$  denote the component group of the correponding  $K_c$ -orbit in  $\mathfrak{s}$ .

**Lemma 5.1.** *If  $\mathcal{O}$  is smoothly cuspidal, then  $\overline{A}(\mathcal{O})$  can be identified with the full component group  $A(\mathcal{O})$ .*

For each finite dimensional representation  $F$ , recall the translation functors

$$T_F : X \mapsto P_\lambda[X \otimes F] \quad (5.2)$$

defined in [BV] Lemma 6.3.

**Theorem 5.3.** Assume that  $\mathcal{O}$  is smoothly cuspidal, and let  $\mathcal{O}(\mathbb{R})$  be a real form.

1. The WF-set of any  $\pi \in \Pi({}^L\mathcal{O})$  is the closure of a single orbit. Denote by  $\Pi(\mathcal{O}(\mathbb{R}))$  the set of representations with WF-set  $\mathcal{O}(\mathbb{R})$ , and by  $\mathcal{G}(\mathcal{O}(\mathbb{R}))$  the corresponding (complexified) Grothendieck group.

2. The number of irreducible representations with WF-set  $\mathcal{O}(\mathbb{R})$  is  $|A(\mathcal{O}(\mathbb{R}))|$ .

3.  $A(\mathcal{O}(\mathbb{R}))$  maps onto  $\overline{A}(\mathcal{O})$ .  $\mathcal{G}(\mathcal{O}(\mathbb{R}))$  decomposes under the (simultaneous) action of  $T_F$  in (5.1) into a direct sum of isomorphic subspaces, each of dimension  $|A(\mathcal{O})|$ .

Given  $\pi \in \Pi({}^L\mathcal{O})$ , we can attach to it a pair of (unions of) real nilpotent orbits,  $(WF(\pi), WF(\tilde{\pi}))$ . This is an invariant for the action of  $T_F$ . The theorem is a consequence of the nontrivial computation of the behaviour of this invariant under induction and coinduction.

## 6. Unitarity

The proof of the unitarity of the distinguished representations attached to cuspidal nilpotent orbits proceeds by induction as in the complex case. Recall that the group is  $Sp(2n, \mathbb{R})$ . To conform to [B], we will write  $g(n)$  for the Lie algebra  $sp(2n, \mathbb{R})$  of rank  $n$ .

Recall that  ${}^L\mathcal{O}$  is an even nilpotent and  $\mathcal{O}$  is its dual. Assume that  $\mathcal{O}$  is smoothly cuspidal. Let  $\pi(\mathcal{O})$  be the *special unipotent spherical* representation in the Arthur L-packet attached to  ${}^L\mathcal{O}$ .

**Proposition 6.1.** The WF-set of  $\pi(\mathcal{O})$  is the closure of the unique  $\mathcal{O}(\mathbb{R})$  for which there exists a Levi component of a real parabolic subalgebra  $m$  such that  $\mathcal{O}(\mathbb{R}) \cap m$  is a principal nilpotent in  $m$ .

**Proposition 6.2.** Assume that  $\mathcal{O}$  is induced (not smoothly). Then the spherical representation in the Arthur L-packet is unitary. More precisely, it occurs as derived functor module from a unitary character on a  $\theta$ -stable parabolic subalgebra satisfying the condition of Theorem 7.1(b) in [V2]. Its unitarity can also be seen from the fact that it occurs at the endpoint of a complementary series of an induced from a spherical distinguished representation on a Levi component of a real proper parabolic subalgebra.

*Proof of Theorem 1.1.* Let  $\mathcal{O}$  with Jordan blocks  $1^{r_1} 2^{r_2} \dots (2k)^{r_{2k}}$  be a special cuspidal nilpotent as at the end of Section 4 (so that  $r_i$  are all even). We do an induction on the number of odd sized Jordan blocks. Assume that none of the  $r_{2l-1} = 0$ , for otherwise the representation is already unitary by Proposition 6.2. The case when  $k = 1$ ,  $r_1 = 2$  can be done by direct calculation.

Let  $r = r_{2k}$ . Denote the spherical distinguished representation attached to  $\mathcal{O}$  by  $\pi({}^L\mathcal{O})$ . Embed the Lie algebra  $m = g(n) \times gl(r+1)$  as a Levi component in  $g(n+r+1)$  and consider the unitarily induced representation

$$\text{Ind}_{g(n) \times gl(r+1)}^{g(n+r+1)} [\pi(\mathcal{O}) \otimes \text{Triv}]. \quad (6.4)$$

The  $WF$ -set is the real form of

$$1^{r_1} 2^{r_2} \dots (2k-2)^{r_{2k-2}} (2k-1)^{r_{2k-1}-2} (2k)^2 (2k+2)^{r_{2k}}, \quad (6.5)$$

which is the support of the spherical distinguished representation corresponding to this nilpotent. The induced representation has two irreducible factors. They are both derived functor induced representations from a parabolic subalgebra with Levi factor  $g(n+1) \times gl(r)$ . The representation being induced is  $\pi(\mathcal{O}') \otimes \chi$ , where  $\pi(\mathcal{O}')$  is the spherical distinguished representation corresponding to the nilpotent orbit

$$1^{r_1} 2^{r_2} \dots (2k-2)^{r_{2k-2}} (2k-1)^{r_{2k-1}-2} (2k)^{r_{2k}+2} \quad (6.6)$$

and  $\chi$  is a unitary character so that the conditions of Theorem 7.1 in [V2] are satisfied. By induction (the occurrence of  $r_{2k-1}$  has been reduced by 2!) both factors are unitary. We only have to show that the signature on the lowest K-type of the other factor is the same as for the spherical K-type. If we write the highest weights of K-types as decreasing sequences of integers, this lowest K-type is

$$(1, \underbrace{\dots, 1}_{r+1}, 0, \dots, 0, \underbrace{-1, \dots, -1}_{r+1}). \quad (6.7)$$

The signature in the induced representation in (6.4) comes from K-types in  $\pi(\mathcal{O})$  of the same type with the number of 1's  $\leq r+1$ . Repeating the same argument using  $g(n) \times gl(r+r_{2k-1}+r_{2k-2}+1)$  as in [B] Section 10, completes the proof.  $\square$

## References

- [A1] Arthur, J.: On some problems suggested by the trace formula. (Lecture Notes in Mathematics, vol. 1041). Springer, Berlin Heidelberg New York 1984, pp. 1–50
- [A2] Arthur, J.: Unipotent automorphic representations: conjectures. Astérisque **171–172**, Orbites unipotentes et représentations II. Groupes  $p$ -adiques et réels, 1989, 13–71
- [ABV] Adams J., Barbasch D., Vogan D.: Special unipotent representation for real reductive groups. Preprint in preparation
- [B] Barbasch D.: The unitary dual for complex classical Lie groups. Invent. math. **96** (1989) 103–176
- [BSS] Barbasch D., Sahi S., Speh B.: Degenerate series representations for  $GL(n, \mathbb{R})$  and Fourier analysis. Symposia Mathematica **XXXI** (1990) 45–69
- [BV] Barbasch D., Vogan D.: Unipotent representations of complex semisimple Lie groups. Ann. Math. **121** (1985) 41–110
- [L] Lusztig G.: Characters of reductive groups over a finite field. Ann. Math. Studies **107**, Princeton University Press, Princeton NJ
- [LV] Lusztig G., Vogan D.: Singularities of closures of K-orbits on flag manifolds. Invent. math. **71** (1983) 365–379
- [M] Moeglin C.: Orbites unipotentes et spectre discret non ramifié. Preprint 1990
- [MW] Moeglin C., Waldspurger J.-L.: Le spectre résiduel de  $GL(n)$ . Preprint 1989
- [S] Sekiguchi J.: Remarks on nilpotent orbits of a symmetric pair. J. Math. Soc. Japan **39** (1987) 127–138
- [V1] Vogan D.: Representations of real reductive Lie groups. Birkhäuser, Boston 1981

- [V2] Vogan D.: Unitarizability of certain series of representations. *Ann. Math.* **120** (1984) 141–187
- [V3] Vogan D.: Irreducible characters of semisimple Lie groups IV. *Duke Math. J.* **49** (1982) 943–1073
- [V4] Vogan D.: Associated varieties and unipotent representations. Preprint 1989
- [VZ] Vogan D., Zuckerman G.: Unitary representations with non-zero cohomology. *Comp. Math.* **53** (1984) 51–90



# Eisenstein Cohomology of Arithmetic Groups and Its Applications to Number Theory

Günter Harder

Mathematisches Institut, Universität Bonn  
Wegelerstraße 10, W-5300 Bonn, Fed. Rep. of Germany

In this note I want to discuss some questions concerning the cohomology of arithmetic groups. They concern the structure of the cohomology and especially the Eisenstein cohomology as a module under the Hecke algebra. My exposition here is rather vague and imprecise. But in a series of papers I investigated these questions in various special cases. Sometimes I proved theorems which give an answer to one of the questions in special cases, sometimes I discuss the conjectures in special situations and make them more precise. I will refer to these papers later on.

My main objective here is to show that a good understanding of these questions will have interesting applications in number theory, especially for the theory of special values of  $L$ -functions.

## 1. The General Setup

We start from a reductive group  $G/\mathbb{Q}$ , let  $G^{(1)}/\mathbb{Q}$  be its derived group, let  $Z/\mathbb{Q}$  be its centre. We denote the adele group of  $G$  by  $G(\mathbb{A})$ , and we decompose it into its finite and infinite part

$$G(\mathbb{A}) = G(\mathbb{R}) \times G(\mathbb{A}_f) = G_\infty \times G(\mathbb{A}_f).$$

We choose a closed subgroup  $K_\infty \subset G_\infty$  whose connected component of the identity is of the form

$$K_\infty^{(0)} = K_\infty^{(1)} \times Z(\mathbb{R})^{(0)},$$

where  $K_\infty^{(1)}$  is the connected component of the identity of a maximal compact subgroup of  $G_\infty^{(1)}$ . Moreover we choose a compact open subgroup  $K_f \subset G(\mathbb{A}_f)$  (the level subgroup), and we consider the space

$$G_\infty/K_\infty \times G(\mathbb{A}_f)/K_f = X_\infty \times G(\mathbb{A}_f)/K_f.$$

The space  $X_\infty$  is a disjoint union of symmetric spaces (the centre contributes by a flat euclidean space). The group  $G(\mathbb{Q})$  acts upon this space. We pass to the quotient and get a space

$$G(\mathbb{Q}) \backslash X_\infty \times G(\mathbb{A}_f)/K_f = \mathcal{S}_K^G.$$

The quotient space  $\mathcal{S}_K^G$  is a locally symmetric space. It is a finite disjoint union

$$\bigsqcup_{\underline{g}_f \in G(\mathbb{Q}) \backslash G(\mathbb{A}_f) / K_f} \Gamma^{(\underline{g}_f)} \backslash X_\infty,$$

where  $\underline{g}_f$  runs over a set of representatives in  $G(\mathbb{A}_f)$  of the double coset space and where  $\Gamma^{(\underline{g}_f)}$  is the arithmetic subgroup

$$\Gamma^{(\underline{g}_f)} = G(\mathbb{Q}) \cap \underline{g}_f K_f \underline{g}_f^{-1}.$$

We call this space the *locally symmetric space attached to  $G/\mathbb{Q}$*  (of level  $K_f$ ).

Finally we choose a rational representation

$$\varrho : G \longrightarrow GL(\mathcal{M})$$

which is defined over  $\mathbb{Q}$  or over some number field (if we want it to be absolutely irreducible) or over  $\overline{\mathbb{Q}}$  (see [Ha1], 1.3). This representation provides a sheaf  $\tilde{\mathcal{M}}$  on  $\mathcal{S}_K^G$  by a standard construction.

Our object of interest are the cohomology groups of  $\mathcal{S}_K^G$  with coefficients in this sheaf

$$H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}).$$

These cohomology groups are finite dimensional vector spaces over the field of definition of  $\varrho$  and they have some very important extra features (1.1 and 1.2).

## 1.1 The Borel-Serre Compactification

In general the quotient space  $\mathcal{S}_K^G$  is not compact. There is a natural construction of a compactification  $i : \mathcal{S}_K^G \longrightarrow \overline{\mathcal{S}_K^G}$ , this is the Borel-Serre compactification (see [B-S]). It is known that  $\overline{\mathcal{S}_K^G}$  is a manifold with corners (provided  $K_f$  is small enough). It is stratified by manifolds  $\partial_P \mathcal{S}_K^G$  which are labelled by the conjugacy classes of parabolic subgroups over  $\mathbb{Q}$ . The closure of a stratum  $\partial_P \mathcal{S}_K^G$  consists of the strata  $\partial_Q \mathcal{S}_K^G$  where  $Q$  runs over the conjugacy classes of rational parabolic subgroups  $Q \subset P$ . One knows that the sheaves  $\tilde{\mathcal{M}}$  extend nicely to  $\overline{\mathcal{S}_K^G}$ , the map  $i : \mathcal{S}_K^G \rightarrow \overline{\mathcal{S}_K^G}$  is a homotopy equivalence and we get an isomorphism

$$H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}) \xrightarrow{\sim} H^*(\overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}).$$

We also consider the cohomology with compact supports  $H_c^*(\mathcal{S}_K^G, \tilde{\mathcal{M}})$ , we get a long exact cohomology sequence

$$\rightarrow H_c^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}) \rightarrow H^*(\overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}) \xrightarrow{r} H^*(\partial \overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}) \rightarrow \dots$$

Let  $M/\mathbb{Q}$  be the Levi quotient of the parabolic subgroup  $P/\mathbb{Q}$ , let  $U_P/\mathbb{Q}$  be the unipotent radical of  $P$  and let  $u_P$  be its Lie algebra. Then one knows that the cohomology groups  $H^*(\partial_P \mathcal{S}_K^G, \tilde{\mathcal{M}})$  can be expressed in terms of the cohomology of the locally symmetric space attached to  $M$  with coefficients in the the Lie

algebra cohomology  $H^*(\mathfrak{u}_P, \tilde{\mathcal{M}})$  (which is a module for  $M$ ). We write without further explanation

$$H^*(\partial_P \mathcal{S}_K^G, \tilde{\mathcal{M}}) = \text{Ind}_P^G H^*(\mathcal{S}_{K^M}^M, \widetilde{H^*(\mathfrak{u}_P, \tilde{\mathcal{M}})}) .$$

(See for instance [Ha1], Theorem 1, [Ha5], 1.5–1.7). The stratum  $\partial_P \mathcal{S}_K^G$  has codimension  $d(P) = \text{corank of } P$  this is the increase of the split rank of  $M$  against the split rank of  $G$ . The covering of  $\partial \mathcal{S}_K^G$  by the  $\partial_P \mathcal{S}_K^G$  yields a spectral sequence with  $E_1^{p,q}$ -term

$$\bigoplus_{P : d(P)=p+1} H^q(\partial_P \mathcal{S}_K^G, \tilde{\mathcal{M}}) \Rightarrow H^{p+q}(\partial \mathcal{S}_K^G, \tilde{\mathcal{M}}) . \quad (\text{Ss})$$

The exact sequence above induces a two-step filtration on the cohomology, but since the cohomology of the boundary will also be filtered, we get a many step filtration on the cohomology.

## 1.2 The Hecke Operators

Let  $\mathcal{H}_{K_f} = \mathcal{C}_c(G(\mathbb{A}_f)/\!/K_f)$  be the space of  $\mathbb{Q}$  (or  $\overline{\mathbb{Q}}$ ) valued functions on  $G(\mathbb{A}_f)$  which have compact support and are biinvariant under  $K_f$ . These functions form an algebra under convolution; this is the so called Hecke algebra. It is a restricted tensor product of local Hecke algebras  $\mathcal{H}_{K_p}$  if  $K_f$  is a product of local groups. We can construct an action of this algebra on all the groups

$$H_c^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}), H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}), H^*(\partial \mathcal{S}_K^G, \tilde{\mathcal{M}})$$

which is compatible with the maps in the cohomology sequence. (To see this we pass to the limit for smaller and smaller level subgroups  $K'_f$ , then the cohomology becomes a  $G(\mathbb{A}_f)$ -module and the  $\varphi \in \mathcal{H}_{K_f}$  act by convolution on

$$H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}) = \varinjlim_{K'_f} H^*(\mathcal{S}_{K'_f}^G, \tilde{\mathcal{M}})^{K_f} .$$

The fundamental problem is to understand the structure of these cohomology groups as modules under the Hecke algebra.

**A Digression.** I will try to make a little bit more precise what it means to understand these modules. To do this I have to explain some ideas which go back to Eichler, Shimura, Deligne, Serre, Langlands and others. The picture that I will give is oversimplified and can only be true in simple cases.

We extend our coefficient system to  $\tilde{\mathcal{M}}_{\mathbb{C}} = \tilde{\mathcal{M}} \otimes \mathbb{C}$ . Then we can compute the cohomology groups  $H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})$  by using differential forms. Within this space of differential forms we have the space of cusp forms which allows us to define a subspace

$$H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})$$

in the cohomology (see [Bo] 5.5, [Schw]). To this space we can apply Hilbert space techniques and get an isotypical decomposition

$$H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}}) = \bigoplus_{\pi} H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi]_{\mathbb{C}}$$

where  $H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi] = H_{\pi}^{m(\pi)}$  is the sum of  $m(\pi)$  copies of the irreducible  $\mathcal{H}$ -module  $H_{\pi}$ . Then  $\pi$  will be the finite component of certain cuspidal automorphic representation  $\pi^* = \pi_{\infty} \otimes \pi$ .

Let us assume that this subspace descends to a subspace  $H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}})$  defined over  $\overline{\mathbb{Q}}$  (or even  $\mathbb{Q}$  if  $\mathcal{M}$  is a  $\mathbb{Q}$ -vector space), this is true in special cases (see [Ha1], 3.2.5 and [Cl]). Then we get for this subspace

$$H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}) \otimes \overline{\mathbb{Q}} = \bigoplus_{\pi} H_{\text{cusp}}^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi],$$

where the isotypical components are now  $\overline{\mathbb{Q}}$  vector spaces (if  $\mathcal{M}$  is a  $\mathbb{Q}$ -vector space, we will get an action of  $\text{Gal}(\overline{\mathbb{Q}}/\mathbb{Q})$  on the set of the contributing  $\pi$ , it will permute the summands.)

There is some general belief that one can attach an arithmetical object  $\mathbf{M}(\pi)$  (something like a motive or a motive with coefficients) to such an isotypical component.

It is not quite clear what such a thing is, but it should have different cohomological realizations, especially it should have  $\lambda$ -adic cohomology groups

$$H^*(\mathbf{M}(\pi) \times \overline{\mathbb{Q}}, \overline{\mathbb{Q}}_{\lambda})$$

which will be modules for  $\text{Gal}(\overline{\mathbb{Q}}/E)$ , where  $E$  is some specific number field.

Then one expects some kind of reciprocity law: For almost all places  $p$  the local module  $H_{\pi_p}$  under the local component  $\mathcal{H}_p$  of the Hecke algebra “is strongly related” to the structure of the  $H^*(\mathbf{M}(\pi) \times \overline{\mathbb{Q}}, \overline{\mathbb{Q}}_{\lambda})$  as a module under  $\prod_{\mathfrak{p}|p} \text{Gal}(\overline{\mathbb{Q}}_{\mathfrak{p}}/E_{\mathfrak{p}})$ . Another way of saying this is the following: Langlands attached to any cuspidal automorphic form  $\pi^*$  a series of  $L$ -functions

$$L(\pi^*, r, s)$$

where  $r$  is a parameter (a representation of the dual group) and where  $s$  is a complex variable (see [La]). These are the *automorphic L-functions*.

On the other hand a motive  $\mathbf{M}$  is also a thing which yields an  $L$ -function

$$L(\mathbf{M}, s)$$

these are the *arithmetic L-functions*. One hopes that the above “strong relationship” says that for a suitable  $r_0$  we have

$$L(\mathbf{M}(\pi), s) \sim L(\pi^*, r_0, s)$$

where  $\sim$  means equality in simple cases.

For more precise information I refer to [Ko, Cl]; this problem has been investigated for Shimura varieties. For the classical case of  $GL_2/\mathbb{Q}$  the existence of  $\mathbf{M}(\pi)$  has been proved by Eichler, Shimura, Deligne and Scholl. In some other cases partial results are known ([Wi, Bl-Ro, Ta]). There exists also some

numerical evidence for the truth of such an assertion in the case where  $\mathcal{S}_K^G$  is not a Shimura variety ([Cr, E-G-M, A-P-T]).

Even if this is rather vague I think we should have this in the back of our mind if we say we want to “understand” the structure of the cohomology as a module under the Hecke algebra. I point out that for the case of the multiplicative group over a number field  $k$  i.e. for  $G/\mathbb{Q} = R_{k/\mathbb{Q}}(G_m)$  the above is the content of class field theory.

### 1.3 Integral Cohomology

If we fix the level subgroup  $K_f$ , and if we choose a  $K_f$ -invariant lattice  $\mathcal{M}_0 \subset \mathcal{M}$ , then it is not difficult to define a sheaf  $\tilde{\mathcal{M}}_0$  on  $\mathcal{S}_K^G$ , and we can define the cohomology groups

$$H_c^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_0), H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}_0), H^*(\partial \overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}_0),$$

and we have the same exact sequence as above. Now we can consider the subalgebra  $\mathcal{H}_{\mathbb{Z}}$  of  $\mathbb{Z}$  valued functions in  $\mathcal{H}$  and with some modification we can define an action of  $\mathcal{H}_{\mathbb{Z}}$  on these integral cohomology groups (see [K-P-S], [Ha5], 1.3). It is certainly also of interest to study the integral cohomology groups as modules under the action of  $\mathcal{H}_{\mathbb{Z}}$ .

## 2. The Eisenstein Cohomology

The aim of the Eisenstein Cohomology is to provide an understanding of the maps

$$H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}}) \xrightarrow{r} H^*(\partial \overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}) \xrightarrow{\delta} H_c^{*+1}(\mathcal{S}_K^G, \tilde{\mathcal{M}})$$

as maps for modules under the Hecke algebra. Actually this of course also requires that we understand the cohomology  $H^*(\partial \overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}})$  as a module under the Hecke algebra and hence at least a full understanding of the cohomology for lower dimensional groups.

We pick a parabolic subgroup  $P$  and look at the associate class  $\mathcal{P} = \{P = P_1, P_2, \dots\}$ . (Two parabolic subgroups are associate if and only if their Levi subgroups are conjugate.) Let  $M$  and  $\mathfrak{u}_P$  be as in 1.1. Let us assume that  $H^*(\mathcal{S}_{KM}^M, H^*(\widetilde{\mathfrak{u}_P}, \mathcal{M}))[\pi_M]_{\mathbb{C}}$  is an isotypical contribution to the cuspidal cohomology  $H_{\text{cusp}}^*(\mathcal{S}_{KM}^M, (\widetilde{\mathfrak{u}_P}, \mathcal{M})_{\mathbb{C}})$ . We consider the group  $W^{(M)} = N_G(M)/M$ , this group gives us a collection of conjugate contributions

$$H^*(\mathcal{S}_{KM}^M, H^*(\widetilde{\mathfrak{u}_P}, \mathcal{M}))[\pi_M]_{\mathbb{C}}$$

(see [Ha1], IV, [Ha3], 1.1.2 for examples) and this gives us a contribution

$$\bigoplus_{P_v \in \mathcal{P}} \bigoplus_w \text{Ind}_{P_v}^G H^*(\mathcal{S}_{KM}^M, H^*(\widetilde{\mathfrak{u}_{P_v}}, \mathcal{M}))[\pi_M]_{\mathbb{C}} \subset \bigoplus_{P_v \in \mathcal{P}} H^*(\partial_{P_v} \overline{\mathcal{S}_K^G}, \tilde{\mathcal{M}}_{\mathbb{C}}).$$

We now want to solve two problems:

The first problem is to attach to our orbit  $[\pi_M]$  a  $\mathcal{H}$ -submodule

$$H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M] \subset H^*(\partial \overline{\mathcal{G}}_K^G, \mathcal{M}_{\mathbb{C}})$$

which should be considered as the contribution of  $[\pi_M]$  to the cohomology of the boundary. This is easily done in the case where  $\mathcal{P}$  consists of maximal parabolic subgroups. In the general case it requires already the solution of the second problem below for groups of smaller semi simple rank. Sometimes this contribution satisfies an extra condition

$$H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M] \quad \begin{array}{l} \text{is a direct summand and does not weakly} \\ \text{intertwine with the complement} \end{array} \quad (\text{bMD})$$

where no weak intertwining means that the two summands do not have any Jordan-Hölder quotient in common.

If the condition (bMD) is satisfied, then it follows easily that the subspace

$$H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M] \subset H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}) \otimes \mathbb{C}$$

is rational, this means that it is defined over  $\overline{\mathbb{Q}}$  and behaves nicely under the action of  $\text{Gal}(\overline{\mathbb{Q}}/\mathbb{Q})$  (see [Ha1], 4.3).

**The second problem** is to understand the intersection

$$\text{Im}(r)(H^*(\mathcal{G}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})) \cap H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M] = H_{\text{glob}}^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M]$$

I will call this the *image of  $H^*(\mathcal{G}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})$  in the subspace  $H^*(\partial \overline{\mathcal{G}}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M]$* .

To attack this second problem we use the theory of Eisenstein series.

As I explained earlier the datum  $[\pi_M]$  yields a collection of spaces

$$\left\{ \text{Ind}_{P_v}^G H_{\text{cusp}}^*(\mathcal{G}_{K_M}^M, H^*(\mathfrak{u}_{P_v}, \widetilde{\mathcal{M}}_{\mathbb{C}})) [w \cdot \pi_M] \right\}_{w \in W^{(M)}, v=1, \dots, s}.$$

Attached to the  $w \cdot \pi_M$  are isotypical spaces  $\mathcal{X}(w \cdot \pi_M)$  of automorphic forms on  $M$  which “induce” certain spaces  $\text{Ind}_{P_v}^G(\mathcal{X}(w \cdot \pi_M)) \subset \mathcal{C}_{\infty}(P_v(\mathbb{Q}) \backslash G(\mathbb{A}) / K_f)$ . The classes in our spaces can be represented by differential forms  $\omega_{\pi, P_v} \in \text{Hom}_{K_{\infty}}(\Lambda^*(\mathfrak{g}/\mathfrak{k}), \text{Ind}_{P_v}^G(\mathcal{X}(w \cdot \pi_M)) \otimes \mathcal{M}_{\mathbb{C}})$ . These forms  $\omega_{\pi, P_v}$  are invariant under  $P_v(\mathbb{Q})$ . By a process of infinite summation over  $P_v(\mathbb{Q}) \backslash G(\mathbb{Q})$  we can try to make them invariant under  $G(\mathbb{Q})$ . But this summation may diverge. We introduce a complex parameter  $\Lambda \in \mathbb{C}^{d(P)}$  and multiply our functions in  $\text{Ind}_{P_v}^G(\mathcal{X}(w \cdot \pi_M))$  by  $\delta_{P_v}^{\Lambda}$ . Then our summation converges for  $\text{Re}(\Lambda)$  in a certain positive cone and defines a holomorphic function in  $\Lambda$  which extends to a meromorphic function for all  $\Lambda$  (see [H-C]). These functions suitably evaluated at  $\Lambda = 0$  provide a space of automorphic forms

$$\text{Eis}^*[\pi_M] \subset \mathcal{A}(G(\mathbb{Q}) \backslash G(\mathbb{A}) / K_f)$$

and a space (see [Ha3], 1.2.1)

$$\text{Eis}^*[\pi_M]_{\mathbb{C}} = \text{Im}(H^*(\mathfrak{g}, K, \text{Eis}^*[\pi_M] \otimes \mathcal{M}_{\mathbb{C}}) \longrightarrow H^*(\mathcal{G}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})).$$

We state the assertion

$$r : \text{Eis}^*[\pi_M]_{\mathbb{C}} \longrightarrow H_{\text{glob}}^*(\partial \mathcal{G}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})[\pi_M] \quad \text{is a surjective map} \quad (\text{Hope})$$

The general philosophy expressed by (Hope) is that the entire cohomology  $H^*(\mathcal{G}_K^G, \tilde{\mathcal{M}}_{\mathbb{C}})$  is built up out of the cuspidal cohomology and the  $\text{Eis}[\pi_M]_{\mathbb{C}}$ .

### 3. Arithmetical Applications

The structure of the module  $\text{Eis}[\pi_M]_{\mathbb{C}}$  as a module under the Hecke algebra will depend in a subtle way on the individual  $\pi_M$ . To be a little more precise: I explained earlier (Section 1) that we can attach various  $L$ -functions  $L(\pi_M^*, r, s)$  to the automorphic form  $\pi_M^*$ . Certain products of some of these  $L$ -functions will occur in the computation of the constant term of the Eisenstein series (see [La]).

The module  $\text{Eis}[\pi_M]_{\mathbb{C}}$  resulting from the above construction will depend on the behavior of the occurring  $L$ -functions at  $s = 0$  (vanishing, poles, special values, for a first subtle case see [Ha3], 3.3). It is my main goal to understand this behavior.

The same also applies to the contribution  $H^*(\partial \mathcal{G}_{K_f}^G, \mathcal{M})[\pi_M]$ , for instance it will be interesting to investigate the influence of the  $L$ -values on the (higher) differentials in (Ss).

I want to explain how such an understanding may have arithmetical consequences.

#### 3.1 Special Values of $L$ -Functions

We assume that our  $[\pi_M]$  satisfies condition (bMD). As I explained earlier we know in this case that

$$H^*(\partial \mathcal{G}_K^G, \tilde{\mathcal{M}})[\pi_M]$$

is defined over  $\overline{\mathbb{Q}}$ , and it is even a rational subspace in the sense of [Ha1], 1.3. Let us assume we have proved that (Hope) is true, then we know that the space

$$r(\text{Eis}^*[\pi_M]_{\mathbb{C}}) \subset H^*(\partial \mathcal{G}_K^G, \tilde{\mathcal{M}})[\pi_M] \otimes \mathbb{C}$$

is a rational subspace. But this subspace can be described in terms of ratios of special values of  $L$ -functions and hence we get rationality results for these values.

The typical result under good circumstances will be that

$$\frac{L(\pi_M^*, r, v - 1)}{L(\pi_M^*, r, v)} \times \pi \times \text{discriminant factor} \in \overline{\mathbb{Q}}^*$$

where  $v, v - 1$  are critical values in the sense of Deligne [De]. This gives us a tool to reduce the proof of Deligne's conjecture to the extreme critical values (for

special examples see [Ha1], Theorem 2, (3), Corollary 4.2.2 and [Ha3], Corollary 3.5.1, it is clear that there will be many more examples where this works).

**3.1.1 The Manin-Drinfeld Principle.** Let us now fix the degree  $v$  of the cohomology, for a given  $[\pi_M]$  we assume (bMD) and (Hope). If the map

$$\mathrm{Eis}^v[\pi_M]_{\mathbb{C}} \xrightarrow{\sim} \mathrm{Im}(r(\mathrm{Eis}^v[\pi_M]_{\mathbb{C}})) \subset H^*(\partial\overline{\mathcal{S}}_K^G; \mathcal{M})[\pi_M] \otimes \mathbb{C},$$

is an isomorphism and if  $\mathrm{Eis}^v[\pi_M]_{\mathbb{C}}$  does not weakly intertwine with its complement, then we say that the *Manin-Drinfeld Principle* (MD) holds (In the classical situations it just follows from the fact that the eigenvalues of the Hecke operators on cusp forms are different from the eigenvalues of Eisenstein classes (multiplicity one)). (The assertion (bMD) is something like (MD) for the cohomology of the boundary).

If we know (MD) then we can conclude that the space  $\mathrm{Eis}^v[\pi_M]_{\mathbb{C}}$ , which is constructed by transcendental means, descends to a subspace defined over  $\overline{\mathbb{Q}}$  and it is actually rational. (See [Ha1], Corollary (4.2.1), (note the subtle point (c)), and [Ha3], Theorem II).

This is an interesting fact by itself, but in some cases we may evaluate these rational classes on certain cycles (modular symbols constructed from subgroups) and if we are lucky we may express the result in terms of special  $L$ -values. Hence we get rationality results for these  $L$ -values (see [Ha1], 5.7.1). As we explained in [H-S] this implies – combined with the results of Don Blasius [Bl] – the truth of Deligne's conjecture [De] in the case of algebraic Hecke characters. (See also [Ha3]).

**3.1.2 Integrality.** We go one step further. Let us assume for the given  $[\pi_M]$  that (bMD) is true, then we can say: If we invert a finite number  $S = \{p_1, p_2, \dots, p_s\}$  of rational primes then we even get a decomposition

$$H^*(\partial\overline{\mathcal{S}}_K^G, \tilde{\mathcal{M}}_{\mathcal{O}_S}) = \text{complement} \oplus H^*(\partial\overline{\mathcal{S}}_K^G, \tilde{\mathcal{M}}_{\mathcal{O}_S})[\pi_M].$$

We assume moreover that (in a fixed degree  $v$ ) (MD) holds. Then we may try to define an isotypical subspace

$$\mathrm{Eis}^v[\pi_M]_{\mathcal{O}_S} \subset H^v(\mathcal{S}_K^G, \mathcal{M}_{\mathcal{O}_S})$$

(There are some problems with torsion which I cannot discuss here), and we may ask for the image

$$r : \mathrm{Eis}^v[\pi_M]_{\mathcal{O}_S} \longrightarrow H_{\mathrm{glob}}^v(\partial\overline{\mathcal{S}}_K^G, \tilde{\mathcal{M}}_{\mathcal{O}_S})[\pi_M].$$

If we have (Hope) then the cokernel of this map will be a finitely generated  $\mathcal{O}_{\mathcal{S}}$ -torsion module. It will be interesting to ask for the structure of this module. It will be nontrivial in many cases, this indicates a failure of an integral version of (MD). In general I hope that the structure of this cokernel may be related to the arithmetic of certain special values of  $L$ -functions of the form  $L(\pi_M^*, r, v)$ . (The conjectures of Deligne on special values say that these values divided by

a suitable period  $\Omega$  are rational numbers. But by construction  $\Omega$  itself is only defined modulo  $\mathbb{Q}^*$ . To give sense to an assertion of the form above one has to define the periods modulo  $\mathcal{O}_S^*$  which is already a problem in itself. Moreover we have to assume that for  $\pi_M$  the problem discussed in *Digression* has been settled.)

An example is discussed [Ha6] in detail: There I take for  $G = GL_2/\mathbb{Q}$ ,  $K_f = GL_2(\hat{\mathbb{Z}})$ , the coefficient system is obtained from the module  $\mathcal{M}_{n,\mathbb{Z}}$  of homogenous polynomials in two variables of degree  $n$  ( $n$  even) with coefficients in  $\mathbb{Z}$ . We invert the primes  $p$  dividing  $n+2$  and those for which  $p-1|n+2$ . Then the structure of the above module is given by a value of the Riemann Zeta function: it is equal to  $\mathbb{Z}_S/\zeta(-1-n)\mathbb{Z}_S$ . (The value of the Zeta function is integral in  $\mathbb{Z}_S$ ). A generalization of this result to a ramified situation will be discussed in the Bonn Diplom-thesis of Ch. Kaiser.

In [Ha2], IV, [Ha4], II, Beispiel  $PGSp_2$  I discuss some other situations where the investigation of this question is of interest.

If our cokernel turns out to be cyclic, i. e. it is isomorphic to  $\mathcal{O}_S/a(\pi_M)\mathcal{O}_S$ , then we can interprete  $a(\pi_M)$  as *the denominator of the Eisenstein class*. This number is of great interest for several reasons. If we pick a class  $\omega \in H_{\text{glob}}^v(\partial \mathcal{P}_K^G, \tilde{\mathcal{M}}_{\mathcal{O}_S})[\pi_M]$  and we look at

$$\text{Eis}(\omega) \in \text{Eis}^v[\pi_M] \subset H^v(\mathcal{P}_K^G, \tilde{\mathcal{M}}),$$

then  $a(\pi_M)\text{Eis}(\omega)$  will be integral. This will give integrality results for special values of  $L$ -functions for instance those discussed in 3.1.1. In our example above this denominator is  $\zeta(-1-n)$  we get integrality results for  $L$ -values of Dirichlet series over real quadratic fields which are well known. The analogous problem for  $GL_2/F$  for an imaginary quadratic extension of  $\mathbb{Q}$  will be discussed in the Bonn dissertation of H. König.

### 3.2 The Mixed Motives

This is even more speculation. As I said already the structure of the space  $\text{Eis}[\pi_M]_{\mathbb{C}}$  will depend on the behavior of certain  $L$ -functions  $L(\pi_M^*, r, s)$  at  $s=0$ , this argument  $s=0$  is the central point (for the functional equation). In [Ha3] I showed that in some cases ( $GL_3/F$ ,  $F$  totally imaginary) this dependence may be quite subtle. The structure of  $\text{Eis}[\pi_M]_{\mathbb{C}}$  depends on whether these  $L$ -values vanish or not and on the sign of the functional equation. I express the hope that for more complicated groups also the order of vanishing will influence the structure of the module (see [Ha3], 3.5.3). It would be exciting if we could read off the order of vanishing of  $L$ -functions in their central point from the structure of certain cohomology groups.

On the other hand one has the Beilinson-Deligne conjectures (originating from the Birch and Swinnerton-Dyer conjectures) which say that these special  $L$ -values should contain arithmetical information (see [R-S-S, Sch]). The order of vanishing should predict the rank of certain Ext-groups in the category of mixed motives. The arithmetic of the values (see remark above) should predict the order or structure of certain finite groups ( $K$ -groups, ideal class groups).

There is some hope to get information of this kind if we take the following detour:

We investigate the structure of  $\text{Eis}[\pi_M]_{\mathbb{C}}$  in dependence of the special value and let the module tell us something about the arithmetic.

Basically the following principle should be adopted in first approximation:

*The vanishing of certain values  $L(\pi_M^*, r, 0)$  may create a failure of the Manin-Drinfeld-principle for  $\text{Eis}[\pi_M]_{\mathbb{C}}$*

This failure has the following effect. To our given  $\pi_M$  there should exist a module  $\mathcal{E}[\pi_M] \subset H^*(\mathcal{S}_K^G, \tilde{\mathcal{M}})$  which is filtered (see 1.1), which is a direct summand in the cohomology and for which no non trivial subquotient occurs in the Jordan-Hölder series of its complement and where finally the top quotient of the filtration maps isomorphically to  $H_{\text{glob}}^*(\partial \mathcal{S}_K^G, \tilde{\mathcal{M}})[\pi_M]$ .

Now we try again to attach a motive  $\mathbf{M}(\mathcal{E}[\pi_M])$  to this module as we did it in our *Digression* (for instance if  $\mathcal{S}_K^G$  is a Shimura variety). But now this motive is also filtered because it inherits the filtration of the cohomology. In contrast to the cuspidal motive which is a pure motive this motive will be an extension of pure motives and we have (hopefully) constructed a mixed motive whose origin lies in the vanishing of an  $L$ -value.

This is discussed in greater detail and for special cases in [Ha4]. For instance in Section II, Beispiel  $PGSp_2$  I show that such a failure of the Manin-Drinfeld principle is created by the Saito-Kurakawa lifting (see [PS]): We start with a modular cusp form for  $Sl_2(\mathbb{Z})$  of weight  $2\text{mod}4$  which is an eigenform for the Hecke operators. In this case we see a mixed motive which is an extension of the Tate motive  $\mathbb{Q}(-2)$  (which is obtained from the contribution of our modular form to the cohomology of the boundary) by a motive attached to the Saito-Kurakawa lifting of our form. This lifting contributes to the cuspidal cohomology of the group  $PGSp_2$ . It should be isomorphic to the motive attached to the original modular form (if we believe the Tate conjecture). In this example I am not able to compute any kind of extension class, i.e. I am not able to decide under which conditions this extension is nontrivial (for a further discussion see [Ha4]) but I can check certain necessary consistencies, especially it is clear the the extension is predicted by the Beilinson-Deligne conjectures.

There is another such construction of a mixed motive for Hilbert modular surfaces (where I can prove nontriviality of the extension). I hope to include into a revised version of [Ha4].

In some sense the above principle has an integral analogue: If certain  $L$ -values are divisible by a prime  $p$  or a power of this prime then the Manin-Drinfeld principle should fail modulo  $p$  or a certain power  $p^\delta$  of  $p$ . This failure is related to the fact that the Eisenstein class picks up a denominator (see 3.1.2). If we multiply it by its denominator to make it integral its restriction to the boundary becomes zero modulo the denominator and hence in a rather imprecise sense it becomes cuspidal modulo the denominator. This gives us congruences between Eisenstein classes and cuspidal cohomology classes which in the classical case of  $Sl_2(\mathbb{Z})$  are related to the classical congruences.

Let us assume that  $p^\delta$  is the exact divisor of  $a(\pi_M)$ . Then we get representations of Galois groups mod  $p^\delta$  which contains the above class  $a(\pi_M)\text{Eis}(\omega)\text{mod } p^\delta$  as an invariant vector by construction but which is bigger. Hence it is an extension of a rank one representation by another representation. This representation

has controlled ramification (this requires deep results from the theory of  $p$ -adic representations of Galois groups) and hopefully it does not split. I discuss this in detail for the values  $\zeta(-1 - n)$  (for even positive  $n$ ) in [Ha6], Chap VI. We get a different approach to some of the results of Mazur-Wiles [M-W] and of Ribet in [Ri].

## References

- [ASL] Automorphic forms, Shimura varieties and  $L$ -functions, I, II. (L. Clozel and J. S. Milne, eds.) Perspectives in Mathematics **10** (1990) 11
- [MLG] Manifolds and Lie Groups. (Hano, J. et al., eds.) Progress in Math. **14** (1981)
- [A-P-T] Ash, A., Pinc, R., Taylor, R.: An  $\hat{A}_4$  extension of  $\mathbb{Q}$  attached to a non-selfdual automorphic form on  $GL(3)$ . Preprint
- [Bo] Borel, A.: Stable real cohomology of arithmetic groups, II. In: [MLG], 21–55
- [Bl] Blasius, Don: On the critical values of Hecke  $L$ -series. Ann. Math. **124** (1986) 23–63
- [Bl-Ro] Blasius, D., Rogawski, J.: Galois representations for Hilbert modular forms. Bulletin A.M.S., July 1989
- [B-S] Borel, A., Serre, J.-P.: Corners and arithmetic groups. Comment. Math. Helv. **48** (1973) 436–491
- [Cl] Clozel, L.: Motifs et formes automorphes. In: [ASL] I, 77–159
- [Cr] Cremona, J. E.: Hyperbolic tessellations, modular symbols, and elliptic curves over complex quadratic fields. Comp. Math. **51** (1984) 275–323
- [De] Deligne, P.: Valeurs de fonctions  $L$  et périodes d'intégrales. Proc. Symp. Pure Math. **33** (part. 2) (1979) 313–346
- [E-G-M] Elstrodt, J., Grunewald, F., Mennicke, J.:  $PSL(2)$  over imaginary quadratic fields. Proc. of the Journée Arithmetique (Metz) (1982), Asterisque
- [H-C] Harish-Chandra Automorphic forms on semisimple Lie groups. (Lecture Notes in Mathematics, vol. 62). Springer, Berlin Heidelberg New York 1968
- [Ha1] Harder, G.: Eisenstein cohomology of arithmetic groups. The case  $GL_2$ . Inv. math. **89** (1987) 37–118
- [Ha2] Harder, G.: Eisensteinkohomologie für Gruppen vom Typ  $GU(2, 1)$ . Math. Ann. **278** (1987) 563–592
- [Ha3] Harder, G.: Some results on the Eisenstein cohomology of arithmetic subgroups of  $GL_n$ . In: Cohomology of arithmetic groups. Proceedings of a conference held at CIRM (J.-P. Labesse, J. Schwermer, eds.). (Lecture Notes in Mathematics). Springer, Berlin Heidelberg New York (to appear)
- [Ha4] Harder, G.: Arithmetische Eigenschaften von Eisensteinklassen, die modulare Konstruktion von gemischten Motiven und von Erweiterungen endlicher Galoismoduln. Preprint, Bonn 1989
- [Ha5] Harder, G.: Eisenstein-Kohomologie arithmetischer Gruppen: Allgemeine Aspekte. Preprint
- [Ha6] Harder, G.: Kohomologie arithmetischer Gruppen. (Textbook in preparation), Chap. VI (preprint)
- [H-S] Harder, G., Schappacher, N.: Special values of Hecke  $L$ -functions and abelian integrals. (Lecture Notes in Mathematics, vol. 1111). Springer, Berlin Heidelberg New York 1985, pp. 17–49
- [Ko] Kottwitz, R.: Shimura varieties and  $\lambda$ -adic representations. [ASL] I, 161–209
- [K-P-S] Kuga, M., Parry, W., Sah: Group Cohomology and Hecke Operators. In: [MLG], 223–266

- [La] Langlands, R.: Euler products. James K. Whittemore lectures, Yale University 1967
- [M-W] Class fields of abelian extensions of  $\mathbb{Q}$ . *Invent. math.* **76** (1984) 179–330
- [PS] Piateski-Shapiro, I.I.: On the Saito-Kurakawa lifting. *Invent. math.* **71** (1983) 309–338
- [R-S-S] Beilinson's conjectures on special values of  $L$ -functions (Rapoport, M., Schappacher, N., Schneider, P., eds.). Perspectives in Mathematics, vol. 4, 1988
- [Ra] Ramakrishnan, D.: Problems arising from the Tate and Beilinson conjectures in the context of Shimura varieties, [ASL] II, 227–252
- [Ri] Ribet, K.: A modular construction of unramified  $p$ -extensions of  $\mathbb{Q}(\mu_p)$ . *Invent. math.* **34** (1976) 151–162
- [Sch] Scholl, A.: Remarks on special values of  $L$ -functions. Preprint
- [Schw] Schwermer, J.: Cohomology of arithmetic groups, automorphic forms and  $L$ -functions. Proceedings of a conference held at CIRM (J.-P. Labesse, J. Schwermer, eds.). (Lecture Notes in Mathematics, vol. 1447). Springer, Berlin Heidelberg New York 1990, pp. 1–29
- [Ta] Taylor, R.: On Galois representations associated to Hilbert modular forms. *Invent. math.* **98** (1989) 265–280
- [Wi] Wiles, A.: On ordinary  $\lambda$ -adic representations associated to modular forms. *Invent. math.* **94** (1988) 503–514

# Crystallizing the $q$ -Analogue of Universal Enveloping Algebras

Masaki Kashiwara

Research Institute for Mathematical Sciences, Kyoto University, Kitashirakawa  
Sakyo-ku, Kyoto 606, Japan

## § 0. Introduction

The notion of the  $q$ -analogue of universal enveloping algebras is introduced independently by Drinfeld and Jimbo in their study of exactly solvable models in statistical mechanics. This algebra  $U_q(\mathfrak{g})$  contains a parameter  $q$  and it becomes the universal enveloping algebra when  $q = 1$ . This parameter is the one of temperature in the context of statistical mechanics and  $q = 0$  corresponds to the absolute temperature zero. Therefore, we can expect that the theory of  $U_q(\mathfrak{g})$  will be simplified at  $q = 0$ . We call the study of  $U_q(\mathfrak{g})$  at  $q = 0$  crystallization. Of course, we cannot deform  $U_q(\mathfrak{g})$  at  $q = 0$ . However, we can construct the bases of representations of  $U_q(\mathfrak{g})$  at “ $q = 0$ ”, and the  $U_q(\mathfrak{g})$ -module structure is described by combinatorics among them. This gives a purely combinatorial description of the tensor category of  $U_q(\mathfrak{g})$ -modules (and hence  $U(\mathfrak{g})$ -modules).

## § 1. Crystal Bases

### 1.1 Definition of $U_q(\mathfrak{g})$

Let us consider the following data:

- (1.1) a finite-dimensional  $\mathbf{Q}$ -vector space  $\mathfrak{t}$ ,
- (1.2) an index set  $I$ ,
- (1.3) a linearly independent subset  $\{\alpha_i; i \in I\}$  of  $\mathfrak{t}^*$  and a subset  $\{h_i; i \in I\}$  of  $\mathfrak{t}$ ,
- (1.4) an inner product  $( , )$  on  $\mathfrak{t}^*$  and
- (1.5) a lattice  $P$  of  $\mathfrak{t}^*$ .

We assume that they satisfy the following conditions:

- (1.6)  $\{\langle h_i, \alpha_j \rangle\}$  is a generalized Cartan matrix (i.e.  $\langle h_i, \alpha_i \rangle = 2$ ,  $\langle h_i, \alpha_j \rangle \in \mathbf{Z}_{\leq 0}$  for  $i \neq j$  and  $\langle h_i, \alpha_j \rangle = 0 \iff \langle h_j, \alpha_i \rangle = 0$ ),
- (1.7)  $(\alpha_i, \alpha_i) \in \mathbf{Z}_{>0}$ ,
- (1.8)  $\langle h_i, \lambda \rangle = \frac{2(\alpha_i, \lambda)}{(\alpha_i, \alpha_i)}$ ,
- (1.9)  $\alpha_i \in P$  and  $h_i \in P^* = \{h \in \mathfrak{t}; \langle h, P \rangle \subset \mathbf{Z}\}$ .

The  $\mathbf{Q}(q)$ -algebra  $U_q(\mathfrak{g})$  is then the algebra generated by the symbols  $e_i, f_i$  ( $i \in I$ ) and  $q^h$  ( $h \in P^*$ ) with the following fundamental relations:

$$(1.10) \quad q^h = 1 \text{ for } h = 0 \text{ and } q^{h+h'} = q^h q^{h'},$$

$$(1.11) \quad q^h e_i q^{-h} = q^{\langle h, \alpha_i \rangle} e_i \text{ and } q^h f_i q^{-h} = q^{-\langle h, \alpha_i \rangle} f_i,$$

$$(1.12) \quad [e_i, f_j] = \delta_{ij} (t_i - t_i^{-1}) / (q_i - q_i^{-1}) \text{ where } q_i = q^{\langle \alpha_i, \alpha_i \rangle} \text{ and } t_i = q^{\langle \alpha_i, \alpha_i \rangle h_i},$$

$$(1.13) \quad \sum_n (-1)^n e_i^{(n)} e_j e_i^{(b-n)} = 0 \text{ and } \sum_n (-1)^n f_i^{(n)} f_j f_i^{(b-n)} = 0 \text{ for } i \neq j \text{ and } b = 1 - \langle h_i, \alpha_j \rangle.$$

Here we used the notations  $[n]_i = (q_i^n - q_i^{-n}) / (q_i - q_i^{-1})$ ,  $[n]_i! = \prod_{k=1}^n [k]_i$  and  $e_i^{(n)} = e_i^n / [n]_i!$ ,  $f_i^{(n)} = f_i^n / [n]_i!$ . We understand  $e_i^{(n)} = f_i^{(n)} = 0$  for  $n < 0$ .

## 1.2 Operators $\tilde{e}_i$ and $\tilde{f}_i$

For a  $U_q(\mathfrak{g})$ -module  $M$  and  $\lambda \in P$ , we set  $M_\lambda = \{u \in M; t_i u = q^{\langle h_i, \lambda \rangle} u\}$  and call it the weight space of weight  $\lambda$ . We say that  $M$  is *integrable* if  $M = \bigoplus_{\lambda \in P} M_\lambda$  and if  $M$  is a union of finitely dimensional sub- $U_q(\mathfrak{g}_i)$ -modules for any  $i$ . Here  $U_q(\mathfrak{g}_i)$  is the subalgebra of  $U_q(\mathfrak{g})$  generated by  $e_i, f_i$  and  $t_i$ .

By the representation theory of  $U_q(sl_2)$ , any element  $u$  of  $M_\lambda$  is uniquely written in the form

$$(1.14) \quad u = \sum f_i^{(n)} u_n \text{ (resp. } = \sum e_i^{(n)} v_n \text{) where } u_n \in \ker e_i \cap M_{\lambda+n\alpha_i} \text{ and } u_n = 0 \text{ except when } n + \langle h_i, \lambda \rangle \geq 0 \text{ and } n \geq 0 \text{ (resp. } v_n \in \ker f_i \cap M_{\lambda-n\alpha_i} \text{ and } v_n = 0 \text{ except when } n \geq \langle h_i, \lambda \rangle \text{ and } n \geq 0).$$

We define the endomorphisms  $\tilde{e}_i$  and  $\tilde{f}_i$  on  $M$  by

$$\tilde{e}_i u = \sum f_i^{(n-1)} u_n$$

and

$$\tilde{f}_i u = \sum f_i^{(n+1)} u_n.$$

Then  $\tilde{e}_i$  and  $\tilde{f}_i$  satisfy the relations symmetric to this:  $\tilde{e}_i u = \sum e_i^{(n+1)} v_n$  and  $\tilde{f}_i u = \sum e_i^{(n-1)} v_n$ .

## 1.3 Crystal Base

Let  $A$  be the subring of  $\mathbf{Q}(q)$  consisting of the rational functions regular at  $q = 0$ . Let  $M$  be an integrable  $U_q(\mathfrak{g})$ -module.

**Definition 1.1.** A *crystal base* of  $M$  is a pair  $(L, B)$  satisfying the following conditions.

$$(1.15) \quad L \text{ is a free sub-}A\text{-module of } M \text{ such that } M = \mathbf{Q}(q) \otimes_A L.$$

$$(1.16) \quad B \text{ is a base of the } \mathbf{Q}\text{-vector space } L/qL.$$

$$(1.17) \quad L = \bigoplus_{\lambda \in P} L_\lambda \text{ and } B = \sqcup_{\lambda \in P} B_\lambda \text{ where } L_\lambda = L \cap M_\lambda, B_\lambda = B \cap (L_\lambda/qL_\lambda).$$

$$(1.18) \quad \tilde{e}_i L \subset L \text{ and } \tilde{f}_i L \subset L. \text{ Hence } \tilde{e}_i \text{ and } \tilde{f}_i \text{ operate also on } L/qL.$$

$$(1.19) \quad \tilde{e}_i B \subset B \cup \{0\} \text{ and } \tilde{f}_i B \subset B \cup \{0\}.$$

$$(1.20) \quad \text{For } b, b' \in B, b' = \tilde{f}_i b \text{ if and only if } b = \tilde{e}_i b'.$$

For a crystal base  $(L, B)$ , the crystal graph is the oriented colored (by  $i \in I$ ) graph with  $B$  as the set of vertices and  $b \xrightarrow{i} b'$  if  $b' = \tilde{f}_i b$ .

The crystal graph describes completely the action of  $\tilde{e}_i$  and  $\tilde{f}_i$  on  $B \sqcup \{0\}$ .

For  $b \in B$ , we set

$$\varepsilon_i(b) = \max\{k \geq 0; \tilde{e}_i^k b \neq 0\}$$

and

$$\varphi_i(b) = \max\{k \geq 0; \tilde{f}_i^k b \neq 0\}.$$

For  $b \in B_\lambda$ , we have

$$\langle h_i, \lambda \rangle = \varphi_i(b) - \varepsilon_i(b).$$

**Example 1.2.** When  $\mathfrak{g} = sl_2$ ,  $U_q(sl_2)$  is the algebra generated by  $e, f, t, t^{-1}$  with the commutation relation  $tet^{-1} = q^2 e, tft^{-1} = q^{-2} f$  and  $[e, f] = (t - t^{-1})/(q - q^{-1})$ . Then any irreducible  $(l + 1)$ -dimensional representation is isomorphic to  $V_l = \bigoplus_{k=0}^l \mathbf{Q}(q)u_k$  with  $fu_k = [k + 1]u_{k+1}$ ,  $eu_k = [l + 1 - k]u_{k-1}$  and  $tu_k = q^{l-2k}u_k$ . Then  $L = \bigoplus A u_k$  and  $B = \{u_k; 0 \leq k \leq l\} \subset L/qL$ . Then  $(L, B)$  is a crystal base of  $V_l$ . Its crystal graph is

$$u_0 \longrightarrow u_1 \longrightarrow \cdots \longrightarrow u_{l-1} \longrightarrow u_l.$$

#### 1.4 Stability by Tensor Product

Let us define the comultiplication of  $U_q(\mathfrak{g})$  by

$$\begin{aligned} \Delta(q^h) &= q^h \otimes q^h, \\ \Delta(e_i) &= e_i \otimes t_i^{-1} + 1 \otimes e_i, \\ \Delta(f_i) &= f_i \otimes 1 + t_i \otimes f_i. \end{aligned}$$

Then  $U_q(\mathfrak{g})$  has a Hopf algebra structure with  $\Delta$  as a comultiplication. By  $\Delta$ , the tensor product of two  $U_q(\mathfrak{g})$ -modules has a structure of  $U_q(\mathfrak{g})$ -module.

**Theorem 1 (Stability by  $\otimes$ ).** Let  $M_1$  and  $M_2$  be two integrable  $U_q(\mathfrak{g})$ -modules and let  $(L_j, B_j)$  be a crystal base of  $M_j$  ( $j = 1, 2$ ). Set  $L = L_1 \otimes_A L_2$  and  $B = \{b_1 \otimes b_2 \in L/qL; b_j \in B_j\}$ .

- (i) Then  $(L, B)$  is a crystal base of  $M_1 \otimes_{\mathbf{Q}(q)} M_2$ .
- (ii) For  $b_j \in B_j$  ( $j = 1, 2$ ), we have

$$\tilde{f}_i(b_1 \otimes b_2) = \begin{cases} \tilde{f}_i b_1 \otimes b_2 & \text{if } \varphi_i(b_1) > \varepsilon_i(b_2) \\ b_1 \otimes \tilde{f}_i b_2 & \text{if } \varphi_i(b_1) \leq \varepsilon_i(b_2), \end{cases}$$

$$\tilde{e}_i(b_1 \otimes b_2) = \begin{cases} b_1 \otimes \tilde{e}_i b_2 & \text{if } \varphi_i(b_1) < \varepsilon_i(b_2) \\ \tilde{e}_i b_1 \otimes b_2 & \text{if } \varphi_i(b_1) \geq \varepsilon_i(b_2). \end{cases}$$

## 1.5 Existence and Uniqueness

We set  $P_+ = \{\lambda \in P ; \langle h_i, \lambda \rangle \geq 0\}$ . For  $\lambda \in P_+$ , let  $V(\lambda)$  be the irreducible integrable  $U_q(\mathfrak{g})$ -module generated by a vector  $u_\lambda$  of weight  $\lambda$  satisfying  $e_i u_\lambda = 0$ . Then

$$V(\lambda) \cong U_q(\mathfrak{g}) / \left( \sum_i (U_q(\mathfrak{g})e_i + U_q(\mathfrak{g})f_i^{1+\langle h_i, \lambda \rangle}) + \sum_{h \in P^*} U_q(\mathfrak{g})(q^h - q^{\langle h, \lambda \rangle}) \right).$$

Let  $L(\lambda)$  be the smallest sub- $A$ -module of  $M$  such that  $L(\lambda)$  contains  $u_\lambda$  and  $L(\lambda)$  is stable by the  $\tilde{f}_i$ .

Let  $B(\lambda)$  be the subset of  $L(\lambda)/qL(\lambda)$  consisting of the non-zero vectors of the form  $\tilde{f}_{i_1} \cdots \tilde{f}_{i_l} u_\lambda$ .

**Theorem 2 (Existence).**  $(L(\lambda), B(\lambda))$  is a crystal base of  $V(\lambda)$ .

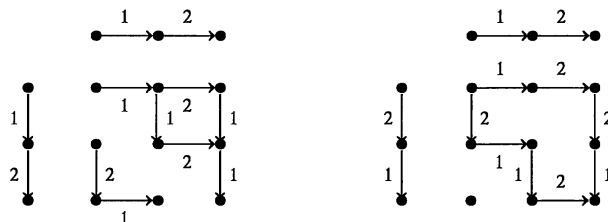
Let  $\mathcal{O}_{\text{int}}$  be the category of integrable  $U_q(\mathfrak{g})$ -modules such that there exists a finite subset  $F$  of  $P$  such that  $M = \bigoplus_{\lambda \in F + Q_-} M_\lambda$ . Here,  $Q_- = \sum \mathbf{Z}_{\leq 0} \alpha_j$ . Then  $\mathcal{O}_{\text{int}}$  is a semi-simple abelian category and any irreducible object is isomorphic to  $V(\lambda)$  for some  $\lambda \in P_+$  ([L], [R]).

**Theorem 3 (Uniqueness).** Let  $(L, B)$  be a crystal base of an object  $M$  in  $\mathcal{O}_{\text{int}}$ . Then there exists an isomorphism  $M \cong \bigoplus_j V(\lambda_j)$  by which  $(L, B)$  is isomorphic to  $\oplus_j (L(\lambda_j), B(\lambda_j))$ .

Combining the Theorems 1, 2 and 3, we can describe completely the tensor category  $\mathcal{O}_{\text{int}}$ .

First note that the crystal graph of  $V(\lambda)$  is connected. Hence the irreducible decomposition of an object in  $\mathcal{O}_{\text{int}}$  is equivalent to the connected component decomposition of the crystal graph. Then Theorem 1 tells us the crystal graph of tensor products

**Example.** Take the case  $\mathfrak{g} = sl_3$  (see §4 for the notation). Let  $\{\Lambda_i\}$  be the dual base of  $\{h_i\}$ . Then the decompositions  $V(\Lambda_1) \otimes V(\Lambda_1) = V(\Lambda_2) \oplus V(2\Lambda_1)$  and  $V(\Lambda_1) \otimes V(\Lambda_2) = V(\Lambda_1 + \Lambda_2) \oplus V(0)$  are described as follows.



## § 2. Crystal Base of $U_q^-(\mathfrak{g})$

### 2.1 Operators $\tilde{e}_i$ and $\tilde{f}_i$ on $U_q^-(\mathfrak{g})$

Let  $U_q^-(\mathfrak{g})$  be the subalgebra of  $U_q(\mathfrak{g})$  generated by the  $f_i$ . Then  $U_q^-(\mathfrak{g})$  has the unique endomorphisms  $e'_i$  and  $e''_i$  such that

$$[e_i, P] = (t_i e''_i(P) - t_i^{-1} e'_i(P))/(q_i - q_i^{-1}) \quad \text{for any } P \in U_q^-(\mathfrak{g}).$$

Then  $e'_i$  and  $f_i$  satisfy the commutation relations:

$$(2.1) \quad e'_i f_j = q_i^{-\langle h_i, \alpha_j \rangle} f_j e'_i + \delta_{ij}.$$

Here we consider  $f_j$  as the left multiplication operator. Then any element  $u$  of  $U_q^-(\mathfrak{g})$  can be uniquely written as

$$u = \sum_{n \geq 0} f_i^{(n)} u_n \quad \text{with} \quad e'_i u_n = 0.$$

We define the endomorphisms  $\tilde{e}_i$  and  $\tilde{f}_i$  of  $U_q^-(\mathfrak{g})$  by

$$\begin{aligned} \tilde{e}_i \left( \sum f_i^{(n)} u_n \right) &= \sum f_i^{(n-1)} u_n \\ \tilde{f}_i \left( \sum f_i^{(n)} u_n \right) &= \sum f_i^{(n+1)} u_n. \end{aligned}$$

Then  $\tilde{e}_i \tilde{f}_i = 1$  holds. Let  $L(\infty)$  be the smallest sub- $A$ -module of  $U_q^-(\mathfrak{g})$  that contains 1 and that is stable by  $\tilde{f}_i$ . Let  $B(\infty)$  be the subset of  $L(\infty)/qL(\infty)$  consisting of the vectors of the form  $\tilde{f}_{i_1} \cdots \tilde{f}_{i_l} \cdot 1$ . Then  $(B(\infty), L(\infty))$  has a similar property to crystal bases.

**Theorem 4.** (i)  $\tilde{e}_i L(\infty) \subset L(\infty)$  and  $\tilde{f}_i L(\infty) \subset L(\infty)$ .

(ii)  $\tilde{e}_i B(\infty) \subset B(\infty) \cup \{0\}$  and  $\tilde{f}_i B(\infty) \subset B(\infty)$ .

(iii)  $B(\infty)$  is a base of  $L(\infty)/qL(\infty)$ .

(iv) If  $b \in B(\infty)$  satisfies  $\tilde{e}_i b \neq 0$ , then  $b = \tilde{f}_i \tilde{e}_i b$ .

The relation of  $(L(\infty), B(\infty))$  and  $(L(\lambda), B(\lambda))$  is given by the following theorem.

**Theorem 5.** For  $\lambda \in P_+$ , let  $\pi_\lambda : U_q^-(\mathfrak{g}) \rightarrow V(\lambda)$  be the  $U_q^-(\mathfrak{g})$ -linear homomorphism sending 1 to  $u_\lambda$ .

(i)  $\pi_\lambda L(\infty) = L(\lambda)$ .

Hence  $\pi_\lambda$  induces the homomorphism  $\bar{\pi}_\lambda : L(\infty)/qL(\infty) \rightarrow L(\lambda)/qL(\lambda)$ .

(ii)  $\{b \in B(\infty); \bar{\pi}_\lambda(b) \neq 0\}$  is isomorphic to  $B(\lambda)$  by  $\bar{\pi}_\lambda$ .

(iii)  $\tilde{f}_i \circ \bar{\pi}_\lambda = \bar{\pi}_\lambda \circ \tilde{f}_i$ .

(iv) For  $b \in B(\infty)$  such that  $\bar{\pi}_\lambda(b) \neq 0$ ,  $\tilde{e}_i \bar{\pi}_\lambda(b) = \bar{\pi}_\lambda(\tilde{e}_i b)$ .

### § 3. Global Crystal Base

Let  $U_q^-(\mathfrak{g})_{\mathbf{Z}}$  be the sub- $\mathbf{Z}[q, q^{-1}]$ -algebra of  $U_q(\mathfrak{g})$  generated by the  $f_i^{(n)}$ . For  $\lambda \in P_+$ , we set  $V_{\mathbf{Z}}(\lambda) = U_q^-(\mathfrak{g})_{\mathbf{Z}} \cdot u_{\lambda}$ . Let  $\bar{\phantom{x}}$  be the ring automorphism of  $U_q^-(\mathfrak{g})$  such that  $\bar{q} = q^{-1}$  and  $\bar{f}_i = f_i$ . This induces the automorphism  $\bar{\phantom{x}}$  of  $V(\lambda)$  by  $\bar{P}u_{\lambda} = \bar{P}u_{\lambda}$  for  $P \in U_q^-(\mathfrak{g})$ .

**Theorem 6.** (i)  $(\mathbf{Q} \otimes U_q^-(\mathfrak{g})_{\mathbf{Z}}) \cap L(\infty) \cap L(\infty)^- \rightarrow L(\infty)/qL(\infty)$  is an isomorphism.  
(ii) For any  $\lambda \in P_+$ ,  $(\mathbf{Q} \otimes V_{\mathbf{Z}}(\lambda)) \cap L(\lambda) \cap L(\lambda)^- \rightarrow L(\lambda)/qL(\lambda)$  is an isomorphism.

Let  $G$  denote the inverse of these isomorphisms. Then we have  $\overline{G(b)} = G(b)$  for  $b \in L(\lambda)/qL(\lambda)$  with  $\lambda \in P_+ \cup \{\infty\}$ . Moreover, we have  $G(b)u_{\lambda} = G(\bar{\pi}_{\lambda}b)$  for any  $b \in L(\infty)/qL(\infty)$ .

**Theorem 7.** For any  $n \geq 0$  and  $i$ ,

$$\begin{aligned} f_i^n U_q^-(\mathfrak{g}) \cap U_q^-(\mathfrak{g})_{\mathbf{Z}} &= \bigoplus_{b \in \tilde{f}_i^n B(\infty)} \mathbf{Z}[q, q^{-1}] G(b), \\ f_i^n V(\lambda) \cap V_{\mathbf{Z}}(\lambda) &= \bigoplus_{b \in \tilde{f}_i^n B(\lambda) \setminus \{0\}} \mathbf{Z}[q, q^{-1}] G(b). \end{aligned}$$

We call  $G(b)$  global crystal base.

It is proven by Lusztig ([L3]) that the canonical bases introduced by himself in [L2] in the case  $A_n$ ,  $D_n$  and  $E_n$  coincide with the global canonical bases introduced here.

### § 4. Example

This example is a joint work with T. Nakashima. Let us take  $\mathfrak{g} = sl_n$ . Hence  $I = \{1, \dots, n-1\}$ ,  $(\alpha_i, \alpha_j) = 1, -1/2, 0$  according to  $i = j$ ,  $|i-j| = 1$ ,  $|i-j| > 0$ . Let  $\Lambda_i \in t^*$  be the dual base of  $h_i$  and take  $\oplus \mathbf{Z}\Lambda_i$  as  $P$ .

Then the crystal graph of  $B(\Lambda_1)$  is

$$\boxed{1} \xrightarrow{1} \boxed{2} \longrightarrow \cdots \longrightarrow \boxed{n-1} \xrightarrow{n-1} \boxed{n} .$$

For  $\lambda = \sum_{v=1}^N \Lambda_{i_v}$  ( $1 \leq i_1 \leq \dots \leq i_N$ ) we embed  $B(\lambda)$  into  $B(\Lambda_1)^{\otimes i_1} \otimes B(\Lambda_1)^{\otimes i_2} \otimes \dots$  by  $u_{\lambda} \mapsto (\boxed{1} \otimes \dots \otimes \boxed{i_1}) \otimes (\boxed{1} \otimes \dots \otimes \boxed{i_2}) \otimes \dots$ . Then  $B(\lambda)$  is parametrized by

$$(\boxed{m_{11}} \otimes \boxed{m_{12}} \otimes \dots \otimes \boxed{m_{1i_1}}) \otimes (\boxed{m_{21}} \otimes \dots \otimes \boxed{m_{2i_2}}) \otimes \dots .$$

in  $B(\Lambda_1)^{\otimes i_1} \otimes \dots$ . We associate to this base the Young diagram  $Y(\lambda)$  with a positive integer in each box as follows

$m_{N1}$		$m_{21}$	$m_{11}$
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	$m_{1i_1}$
		$m_{2i_2}$	
$m_{N,i_N}$			

Here  $Y(\lambda)$  is the Young diagram with the columns with length  $i_1, \dots, i_N$ .

**Theorem.** *By this correspondence,  $B(\lambda)$  is equal to the set of semi-standard tableaux with shape  $Y(\lambda)$  (i.e.  $\{m_{ij}\}$  satisfies  $m_{ij} \leq m_{i'j}$  if  $i > i'$  and  $m_{ij} < m_{i'j}$  if  $j < j'$ ).*

## § 5. Remarks

The notion of crystal base is introduced in [K<sub>1</sub>] under the form dual to the one given here. Theorem 2,3 in the case of  $A_n, B_n, C_n$  and  $D_n$  and Theorem 1 are proven there. In [M], the crystal graph of basic representation of  $U_q(\widehat{\mathfrak{sl}}_n)$  is given. The results here have been announced in [K<sub>2</sub>]. Independently, Lusztig introduced the notion of canonical bases in the case  $A_n, D_n, E_n$  ([L<sub>2</sub>]) and he showed that they coincide with global canonical bases ([L<sub>3</sub>]).

## References

- [D] Drinfel'd, V. G.: Hopf algebra and the Yang-Baxter equation. Sov. Math. Dokl. **32** (1985) 254–258
- [J] Jimbo, M.: A  $q$ -difference analogue of  $U(g)$  and the Yang-Baxter equation. Lett. Math. Phys. **10** (1985) 63–69
- [K] Kashiwara, M.: 1. Crystallizing the  $q$ -analogue of universal enveloping algebras. Commun. Math. Phys. **133** (1990) 249–260  
2. Bases crystallines. C. R. Acad. Sci. Paris **311** (1990) 277–280
- [L] Lusztig, G.: 1. On quantum groups. J. Algebra (1990)  
2. Canonical bases arising from quantized enveloping algebra. J. AMS  
3. Canonical bases arising from quantized enveloping algebra, II. Preprint
- [M] Misra, K. C., Miwa, T.: Crystal bases for basic representation of  $U_q(\widehat{\mathfrak{sl}}(n))$ . Commun. Math. Phys. **134** (1990) 79–88
- [R] Rosso, M.: Analogue de la forme de Killing et du théorème d’Harish-Chandra pour les groupes quantiques. Ann. Sci. Ec. Norm. Sup. **23** (1990) 445–467



# Classification of Simple Graded Lie Algebras of Finite Growth

Olivier Mathieu \*

D. M. I., ENS, 45, rue d'Ulm, F-75005 Paris, France, and  
Rutgers University, Department of Mathematics, Hill Center  
New Brunswick, NJ 08903, USA

## Introduction

*Conventions:* The ground field is  $\mathbb{C}$ . By LA we mean Lie algebra.

Let us start with a few definitions.

- A LA  $\mathcal{L}$  endowed with a decomposition

$$\mathcal{L} = \bigoplus_{n \in \mathbb{Z}} \mathcal{L}_n$$

is called a *graded LA* if we have  $[\mathcal{L}_n, \mathcal{L}_m] \subseteq \mathcal{L}_{n+m}$ . Moreover we will always assume that  $\dim \mathcal{L}_n < \infty$  for any  $n < \infty$ . With our convention any graded LA is an *ordinary* LA and the notion should not be confused with super LA which are often called graded LA as well.

- A subspace  $V$  of  $\mathcal{L}$  is called *homogeneous* if we have

$$V = \bigoplus_{n \in \mathbb{Z}} V_n$$

(where  $V_n = V \cap \mathcal{L}_n$ ). The LA  $\mathcal{L}$  is called *simple graded* if any homogenous ideal is trivial (i.e.  $O$  or  $\mathcal{L}$ ) and if  $\dim \mathcal{L} \geq 2$ .

- Say that  $\mathcal{L}$  has *finite growth* if

$$\dim \mathcal{L}_n \leq P(n)$$

for some polynomial  $P$ .

We have recently proved the following theorem [M2].

**Theorem (1990).** *Let  $\mathcal{L}$  be a simple graded LA of finite growth. Then  $\mathcal{L}$  is isomorphic to one of the following LA:*

1. A simple finite dimensional LA
2. An affine LA

---

\* Work done under the hospitality of IAS at Princeton. I thank IAS for its support (DMS Grant 8610730).

3. A LA of Cartan type
4.  $\mathbf{W}$  (Virasoro-Witt LA).

The previous theorem has been conjectured by V.G. Kac. Alltogether there are 14 infinite series and 13 exceptional LA. In part 1 (zoology) we will give precise definitions of the involved LA. Before we would like to make a few remarks. The origin of Kac conjecture comes from the following result [K].

**Theorem** (V. G. Kac, 1967). *Let  $\mathcal{L}$  be a simple graded LA of finite growth. Assume*

(\*)  $\mathcal{L}$  is generated by its “local part”  $\mathcal{L}_{-1} \oplus \mathcal{L}_0 \oplus \mathcal{L}_1$

(\*\*) the  $\mathcal{L}_0$ -module  $\mathcal{L}_{-1}$  is irreducible.

*Then  $\mathcal{L}$  is isomorphic to*

1. a finite dimensional LA,
2. an affine LA or
3. a Cartan type LA.

Moreover it follows from 1967 Kac paper the existence of “continous families” of simple graded LA. Thus there are *no hopes* for a classification *without the growth hypothesis*. Note also that in characteric  $p \neq 0$  the classification of finite dimensional simple LA is still open.<sup>1</sup>

## Part 1: Zoology

In the section we will describe some species i.e. the LA involved in the Theorem. Altough each of them admit infinitely many different gradings we can describe all of them. For the simplicity of the exposition we will describe theses gradings in one case only.

### (1.1) Finite Dimensional Simple LA

Recall that finite dimensional simple Lie algebras have been classified around 1900 by Killing and Cartan. Four infinite series and five exceptional Lie algebras occur in their classification. The LA of the four infinite series are called classical LA. They are the following one.

- $A_n$  or  $\mathfrak{sl}(n+1)$
- $B_n$  or  $\mathfrak{so}(2n+1)$
- $C_n$  or  $\mathfrak{sp}(2n)$
- $D_n$  or  $\mathfrak{so}(2n)$

It is not easy to give a simple description of the five exceptional simple LA  $E_6, E_7, E_8, F_4$  and  $G_2$

---

<sup>1</sup> At ICM conference G. Seligman tells us that H. Strade and R. Wilson have recently announced the classification of finite dimensional simple LA over field of characteristic  $p > 7$ .

### (1.2) Affine (Kac-Moody) LA

Let  $\mathfrak{g}$  be a finite dimensional simple LA, let  $\omega$  be an automorphism of  $\mathfrak{g}$  of finite order  $\ell$  and let  $\eta$  be  $\ell$ -root of unity. Set  $L(\mathfrak{g}) = \mathfrak{g} \otimes \mathbb{C}[t, t^{-1}]$ . Define the automorphism  $\tilde{\omega}$  of  $L(\mathfrak{g})$  by:

$$\tilde{\omega}(g \otimes t^n) = \eta^n \omega(g) \otimes t^n.$$

Let  $L(\mathfrak{g}, \omega, \eta)$  be the LA of fixed points under  $\tilde{\omega}$ . A LA isomorphic to some  $L(\mathfrak{g}, \omega, \eta)$  is called *affine*. The definition is not accurate because there are many non-trivial isomorphisms between various  $L(\mathfrak{g}, \omega, \eta)$ . Fortunately V.G. Kac found a one to one parametrization of these isomorphism classes [K]. Actually he proved that affine LA are exactly parametrized by automorphisms of Dynkin diagrams. All together there are 6 infinite series and 7 exceptional affine LA. With the usual notation affine algebras are

$$A_n^{(1)}, B_n^{(1)}, C_n^{(1)}, D_n^{(1)}, A_n^{(2)}, D_n^{(2)}, D_{(4)}^{(3)}, E_{(6)}^{(1)}, E_{(6)}^{(2)}, E_{(7)}^{(1)}, E_{(8)}^{(1)}, G_{(2)}^{(1)}, F_{(4)}^{(1)}.$$

Affine LA are also called loop algebras because any element of  $L(\mathfrak{g})$  can be identified with a  $\mathfrak{g}$ -valued map on  $S^1$  whose Fourier decomposition is finite.

### (1.3) Cartan Type LA

Let  $\mathbf{W}_n$  be the LA of derivations of the ring of polynomials  $\mathbb{C}[X_1, \dots, X_n]$ . Thus an element  $\partial$  of  $\mathbf{W}_n$  is a vector field with polynomial coefficients. Note  $\text{Lie}(\partial)$  be the Lie derivative action on spaces of differential forms.

Set  $\mathbf{S}_n = \{\partial \in \mathbf{W}_n | \text{Lie}(\partial) \cdot v = 0\}$  where  $v$  is the usual volume form  $dX_1 \wedge \dots \wedge dX_n$ .

For  $n = 2m$  let  $\omega = \sum_{1 \leq i \leq m} dX_i \wedge dX_{m+i}$  be the usual symplectic form and set  $\mathbf{H}_n = \{\partial \in \mathbf{W}_n | \text{Lie}(\partial) \cdot \omega = 0\}$ .

For  $n = 2m + 1$  let  $\alpha = dX_n + \sum_{1 \leq i \leq m} X_i dX_{m+i} - X_{m+i} dX_i$  be the usual contact 1-form and set  $\mathbf{K}_n = \{\partial \in \mathbf{W}_n | \alpha \wedge \text{Lie}(\partial) \cdot \alpha = 0\}$ .

The LA  $\mathbf{W}_n, \mathbf{S}_n, \mathbf{H}_n, \mathbf{K}_n$  are called *Cartan type LA*. These four infinite series have been discovered by Cartan around 1910 [C].

### (1.4) The Virasoro-Witt LA

Let  $\mathbf{W}$  be the LA of derivations of  $\mathbb{C}[T, T^{-1}]$ .

**Remarks:** 1. Each of the previous LA admits *infinitely many* gradings but only *finitely many* of them satisfy the hypotheses (\*) (\*\*\*) of Kac theorem.

Example for  $\mathcal{L} = \mathbf{W}_n$ . Let  $\underline{a} = a_1, \dots, a_n$  be a sequence of non zero integers of same sign. There is a unique grading of  $\mathcal{L}$  such that the element  $X_1^{m_1} \cdots X_n^{m_n} \partial / \partial X_j$

is homogeneous of degree  $a_1m_1 + \dots + a_nm_n - a_j$ . For any  $\sigma \in S_n$  (where  $S_n$  is the symmetric group) the gradings associated with  $a_1, \dots, a_n$  and with  $a_{\sigma(1)}, \dots, a_{\sigma(n)}$  are obviously isomorphic. It is easy to prove that the induced map from  $(\mathbb{Z}_+^n \sqcup \mathbb{Z}_-^n)/S_n$  to the set of isomorphism classes of gradings of  $\mathcal{L}$  is one to one and onto. But the grading associated with  $(1, \dots, 1)$  is the *only one* satisfying Kac's hypotheses (\*), (\*\*\*) ( $n \geq 2$ ).

*Example:* No grading of  $\mathbf{W}, \mathbf{W}_1$  satisfies (\*).

2. The terminology “affine LA”, “Virasoro LA” is often used for central extension of LA considered here.

3. Affine LA are “simple graded” but not simple (because of evaluation maps  $L(\mathfrak{g}) \rightarrow \mathfrak{g}$ ). Conversely a simple graded LA which is not simple is affine (it is a statement).

**Table I.** Simple graded LA of finite growth

dim $\mathcal{L} < \infty$	$A_n B_n C_n D_n$ $E_6 E_7 E_8 F_4 G_2$
Affine	$A_n^{(1)} B_n^{(1)} C_n^{(1)} D_n^{(1)}$ $A_n^{(2)} D_n^{(2)}$ $D_4^{(3)} E_6^{(1)} E_6^{(2)} E_7^{(1)}$ $E_8^{(1)} F_4^{(1)} G_2^{(1)}$
Cartan type	$\mathbf{W}_n \mathbf{S}_n \mathbf{H}_n \mathbf{K}_n$
Virasoro-Witt	$\mathbf{W}$

## Part 2: About the Proofs

In this section we will describe the general plan of the proof. The proof divides into 3 Steps.

*Step 1:* Define 4 abstract classes of graded LA.

*Step 2: Lemma:* Any simple graded LA belongs to one of the 4 previous classes.

*Step 3:* Show 4 classification theorems (i.e. one for each class).

**Step 1. 4 Definitions:** To meet step 1 we will define four abstract classes of graded Lie algebras.

Let  $h$  be a finite dimensional nilpotent LA and let  $M$  be a finite dimensional  $h$ -module. Recall that  $M$  decomposes as

$$M = \bigoplus M^\lambda$$

where  $\lambda$  runs over  $h^*$  and  $M^\lambda$  is the generalized eigenspace associated with  $\lambda$  (Engel Theorem). Let  $\mathfrak{g}$  be a finite dimensional LA. A *Cartan subalgebra* (or CSA) is a nilpotent subalgebra equal to its normalizer. It is classical that CSA

do exist. Moreover any two CSA are conjugated under a product of elementary automorphisms. Let  $\mathcal{L}$  be a graded LA. Pick a CSA  $h$  of  $\mathcal{L}_0$  and consider each  $\mathcal{L}_n$  as an  $h$ -module. Thus we have:

$$\mathcal{L} = \bigoplus_{\substack{n \in \mathbb{Z} \\ \lambda \in h^*}} \mathcal{L}_n^\lambda.$$

Set:  $\Delta = \{(n, \lambda) / \mathcal{L}_n^\lambda \neq 0\}$ .

Let  $Q$  be the subgroup of  $\mathbb{Z} \times h^*$  generated by  $\Delta$ .

**Definition 1.**  $\mathcal{L}$  is called *without roots* iff  $\Delta \subseteq \mathbb{Z} \times 0$

**Definition 2.**  $\mathcal{L}$  is called *weakly integrable* iff

1)  $\Delta \not\subseteq \mathbb{Z} \times 0$

2)  $\bigcap_{s \geq 1} \text{ad}^s(\mathcal{L}_n^\lambda). \mathcal{L} = 0$

for any  $(n, \lambda) \in \Delta, \lambda \neq 0$ .

Set:  $\mathcal{L}^+ = \bigoplus_{s > 0} \mathcal{L}_s, \mathcal{L}^- = \bigoplus_{s < 0} \mathcal{L}_s$ .

**Definition 3.** Say that  $\mathcal{L}$  is of type  $\mathcal{C}$  iff  $\dim \mathcal{L}^+$  or  $\dim \mathcal{L}^-$  is finite but  $\dim \mathcal{L}$  is infinite.

A subset  $X$  of  $Q$  is called *quasi-order* iff

$\forall \tilde{\alpha} \in Q \exists N \geq 0 \forall m \geq N \forall \tilde{\beta}_1 \dots \tilde{\beta}_m \in X$  we have  $\tilde{\alpha} + \tilde{\beta}_1 + \dots + \tilde{\beta}_m \in X$ .

Let  $\tilde{\alpha} \in Q$ . The LA  $\mathcal{L}$  is called  $\tilde{\alpha}$ -deep if we have  $[\mathcal{L}_X, \mathcal{L}] = \mathcal{L}$  for any quasi-order  $X$  such that  $X \cup \{\tilde{\alpha}\}$  is still a quasi-order (by definition we set  $\mathcal{L}_X = \bigoplus_{(n, \lambda) \in X} \mathcal{L}_n^\lambda$ ).

**Definition 4.**  $\mathcal{L}$  is called *deep* if  $\mathcal{L}$  is  $\tilde{\alpha}$ -deep for some  $\tilde{\alpha} = (n, \lambda) \in \Delta$  with  $n \neq 0, \lambda \neq 0$ .

Thus Definitions 1, 2, 3, 4 define 4 abstract classes of graded LA. Moreover any two CSA of  $\mathcal{L}_0$  are conjugated under a degree 0 automorphism of  $\mathcal{L}$ . Hence the definitions do not depend on a choice for  $h$ .

## Step 2

**Lemma 1.** Any simple graded LA  $\mathcal{L}$  satisfies *exactly one* of the following assertions.

- 1)  $\mathcal{L}$  is without roots.
- 2)  $\mathcal{L}$  is weakly integrable.
- 3)  $\mathcal{L}$  is of type  $\mathcal{C}$ .
- 4)  $\mathcal{L}$  is deep.

### Step 3. 4 Classification Theorems

The previous lemma splits the category of simple graded LA into four subcategories. Each of the following four theorem is a classification theorem for each of the four classes. Thus The main theorem is an obvious corollary of these four theorems.

**Theorem 1.** *Let  $\mathcal{L}$  be a simple graded LA without roots. Then  $\mathcal{L}$  has infinite growth.*

**Theorem 2.** *Let  $\mathcal{L}$  be a weakly integrable simple graded LA. Then  $\dim \mathcal{L} < \infty$  or  $\mathcal{L}$  is affine.*

**Theorem 3.** *Let  $\mathcal{L}$  be a simple graded LA of type C. The  $\mathcal{L}$  is of Cartan type.*

**Theorem 4.** *Let  $\mathcal{L}$  be a deep simple graded LA of finite growth. Then  $\mathcal{L}$  is isomorphic to  $W$ .*

Theorem 1 has the following consequence. Any simple graded LA  $\mathcal{L}$  with  $\mathcal{L}_0 = 0$  has infinite growth. The growth hypothesis in theorem 1 is crucial because there are simple graded LA  $\mathcal{L}$  with  $\mathcal{L}_0 = 0$ . However there are no growth hypotheses for Theorems 2, 3.

Thus Theorems 1, 2, 3, 4 and the Lemma implies the classification of simple graded LA of finite growth.

#### Some References:

Say that  $\mathcal{L}$  has *growth*  $\leq 1$  if we have  $\dim \mathcal{L}_n \leq C$  for some constant  $C$ . In a previous paper [M1] we classify simple graded LA of growth  $\leq 1$ .

- 1) The proof of Theorem 1 follows the same line as Theorem 1 in [M1].
- 2) 90 % of Theorem 2 was already proved in [M1].
- 3) The proof of Theorem 3 essentially uses homological Kostant formula, Kac theorem and a calculation of characteristic variety (following a nice trick of V. Guillemin).
- 4) The proof of Theorem 4 is the main difficulty. At some point we use Gabber-Kac theorem. Otherwise it is elementary.

## Part 3: More About the Proofs

The proof of the theorem is quite long. The main “tools” are the following ones:

- 1) *Basic Tool:* We get informations from any “formal construction” of ideals. Obvious examples are centers, derived algebras... We can also use the notion of “quasi-order” for that purpose. Another typical example is the following. Let  $\mathcal{L}$  be any graded LA.

**Lemma 2.** Assume that  $\mathcal{L} = \mathcal{A} + \mathcal{B}$  where  $[\mathcal{A}, \mathcal{A}] \subseteq \mathcal{A}$  and  $[\mathcal{A}, \mathcal{B}] \subseteq \mathcal{B}$ . Then  $\mathcal{B} + [\mathcal{B}, \mathcal{B}]$  is an ideal.

The lemma is obvious but it is used many time to construct subalgebras which behaves like  $sl(2)$  or like an Heisenberg algebra.

2) *Another Tool: Partial LA.* Let  $a < 0$ ,  $b > 0$  be integers. A *partial LA* is a graded vector space

$$\Gamma = \bigoplus_{a \leq i \leq b} \Gamma_i$$

endowed with *partial brackets*  $\Gamma_i \times \Gamma_j \rightarrow \Gamma_{i+j}$  (for  $a \leq i, j, i + j \leq b$ ) satisfying partial Jacobi identities. For a graded LA  $\mathcal{L}$ , its *partial part*

$$\text{Part } \mathcal{L} = \bigoplus_{a \leq i \leq b} \mathcal{L}_i$$

is a partial LA. Conversely any partial LA  $\Gamma$  is the partial part of some graded LA  $\mathcal{L}$ . Among such  $\mathcal{L}$ 's there is a *minimal model*  $\mathcal{L}_{\min}(\Gamma)$ .

**Lemma 3.** Let  $\mathcal{L}$  be a graded LA and  $\Gamma$  be a partial LA. If  $\Gamma$  is a subquotient of Part  $\mathcal{L}$  then  $\mathcal{L}_{\min}(\Gamma)$  is a subquotient of  $\mathcal{L}$ .

Especially if  $\mathcal{L}_{\min}(\Gamma)$  does have infinite growth,  $\mathcal{L}$  does. It allows us to restrict the possible partial part of graded LA of finite growth because we prove that for particular 6 series of partial LA  $\Gamma$  their models  $\mathcal{L}_{\min}(\Gamma)$  have infinite growth.

3) *Another Tool: Ranks.* The *rank* of a graded LA  $\mathcal{L}$  is the dimension over  $\mathbb{Q}$  of  $\mathbb{Q} \otimes_{\mathbb{Z}} \mathcal{L}$ . In the proof of Theorem 4 we study two cases :

$$1) \text{ rank } \mathcal{L} = 1$$

$$2) \text{ rank } \mathcal{L} \geq 2$$

Actually the  $\text{rank } \geq 2$  case is by far easier.

4) *Last Tool: Coadjoint Estimates.* Let  $\mathcal{L}^* = \bigoplus \mathcal{L}_n^*$  be the *graded dual* of the graded LA  $\mathcal{L}$ . For a homogenous  $\zeta \in \mathcal{L}^*$  the space  $\mathcal{L} \cdot \zeta \subseteq \mathcal{L}^*$  is homogenous. For a simple graded LA  $\mathcal{L}$  it is easy to show that the growth of  $\mathcal{L} \cdot \zeta$  is *independant* of  $\zeta \neq 0$ . The following lemma is very crucial in proving Theorem 4.

**Lemma 4.** Assume  $\mathcal{L}$  deep, simple graded and of *finite growth*. Then  $\mathcal{L} \cdot \zeta$  has growth exactly 1.

**Table II.** Simple graded Lie algebras

Type	Finite growth	Infinite growth
Without root	$\emptyset$	A lot [Continuous families]
Weakly integrable	Finite dim. or Affine	$\emptyset$
Type $\mathcal{C}$	Cartan type	$\emptyset$
Deep	Virasoro-Witt	A lot [Continuous families]

## References

- [C] E. Cartan: Les groupes de transformations continus infinis simples. Ann. Sci. Ecole Norm. Sup. **26** (1909) 93–161
- [K] V. G. Kac: Simple graded Lie algebras of finite growth. Math. USSR Izv. **2** (1968) 1271–1311
- [M1] O. Mathieu: Classification des algèbres de Lie graduées simples de croissance  $\leq 1$ . Invent. math. **86** (1986) 371–426
- [M2] O. Mathieu: Classification of simple graded Lie algebras of finite growth. Preprint 1990

# Orbits on Flag Manifolds

Toshihiko Matsuki

College of Liberal Arts and Sciences, Kyoto University, Kyoto 606, Japan

## 1. $H$ -Orbits on $X = G/P$

Let  $G$  be a connected real semisimple Lie group and  $X$  the flag manifold of  $G$ .  $X$  is a homogeneous space of  $G$  and the isotropy subgroup  $P = P_x$  of each point  $x$  of  $X$  is called a minimal parabolic subgroup of  $G$ . Let  $\sigma$  be an involutive automorphism ( $\sigma^2 = id.$ ) of  $G$  and  $H$  a subgroup of  $G^\sigma = \{x \in G \mid \sigma x = x\}$  containing the identity component  $G_0^\sigma$  of  $G^\sigma$ . Irreducible pairs  $(g, h)$  of Lie algebras of  $G$  and  $H$  are classified by [Be].

The following are special cases of  $H$ -orbit decompositions of  $X = G/P$ .

(i) Let  $\sigma$  be a Cartan involution of  $G$ ,  $\mathfrak{g} = \mathfrak{k} \oplus \mathfrak{s}$  the Cartan decomposition of the Lie algebra  $\mathfrak{g}$  of  $G$  for  $\sigma$  and  $K = H = G^\sigma$ . Then  $P = MAN$  where  $A = P \cap \exp \mathfrak{s}$ ,  $M = Z_K(A)$  and  $N$  is the unipotent radical of  $P$ . The Iwasawa decomposition  $G = KAN (\cong K \times A \times N)$  implies that  $\#(K \setminus G/P) = 1$ .

(ii) Let  $G = G_1 \times G_1$ ,  $P = P_1 \times P_1$  and  $\sigma(x, y) = (y, x)$  for  $(x, y) \in G_1 \times G_1$ . Then  $H = G^\sigma = \{(x, x) \in G \mid x \in G_1\}$ . Since  $H \setminus G \cong G_1$  by the map  $H(x, y) \mapsto x^{-1}y$ , the double coset decomposition  $H \setminus G/P$  is identified with the Bruhat decomposition  $P_1 \setminus G_1/P_1$ .

(iii) When  $G$  is a complex semisimple Lie group and  $\sigma$  is a conjugation of  $G$ ,  $H$ -orbits on  $X$  are studied in [A]. This study suggested the formulation for the following general cases.

Let  $\theta$  be a Cartan involution of  $G$  such that  $\sigma\theta = \theta\sigma$ ,  $\mathfrak{g} = \mathfrak{k} \oplus \mathfrak{s}$  the Cartan decomposition of  $\mathfrak{g}$  for  $\theta$  and  $K = G^\theta$ .

**Definition.** An element  $x$  of  $X$  is called “special” when  $A_x = P_x \cap \exp \mathfrak{s}$  is  $\sigma$ -stable. Put

$$U = \{x \in X \mid x \text{ is special}\}.$$

**Theorem 1** [R, M1].  $K \cap H \setminus U \cong H \setminus X$  by the inclusion map  $U \hookrightarrow X$ .

There exists a unique subgroup  $H^a$  of  $G$  such that  $G_0^{\sigma\theta} \subset H^a \subset G^{\sigma\theta}$  and that  $K \cap H^a = K \cap H$ . (Note  $(H^a)^\theta = H$ .)

**Corollary** [M1]. There exists a one-to-one correspondence  $D \mapsto D^a$  between  $H$ -orbits and  $H^a$ -orbits on  $X$  given by  $K \cap H \setminus U \cong H \setminus X$  and  $K \cap H \setminus U \cong H^a \setminus X$ .

*Example 1.* Let  $G = SL(2, \mathbb{C})$ . Then  $X = P^1(\mathbb{C}) = \mathbb{C} \cup \{\infty\}$ ,

$$\text{where } \begin{pmatrix} a & b \\ c & d \end{pmatrix} x = \frac{ax+b}{cx+d} \quad \text{for } \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in SL(2, \mathbb{C}), x \in X.$$

$$\text{Let } \sigma \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} a & -b \\ -c & d \end{pmatrix}, \text{ and } \theta g = {}^t \bar{g}^{-1}.$$

Then

$$K = SU(2), \quad H = G^\sigma = \left\{ \begin{pmatrix} a & 0 \\ 0 & a^{-1} \end{pmatrix} \mid a \in \mathbb{C}^\times \right\},$$

$$H^a = G^{\sigma\theta} = SU(1, 1) = \left\{ \begin{pmatrix} a & b \\ \bar{b} & \bar{a} \end{pmatrix} \mid a\bar{a} - b\bar{b} = 1 \right\}$$

The  $H$ -orbits on  $X$  are  $\{0\}$ ,  $\mathbb{C}^\times$  and  $\{\infty\}$  and the corresponding  $H^a$ -orbits are  $\{|x| < 1\}$ ,  $\{|x| = 1\}$  and  $\{|x| > 1\}$ , respectively. ( $U = \{0\} \cup \{|x| = 1\} \cup \{\infty\}$ .)

## 2. Expression by Symbols

*Remark 1.* If  $H = G_0^\sigma$ , then  $H \setminus X$  depends only on the pair  $(g, \sigma)$  because

$$X \cong \text{the set of minimal parabolic subalgebras of } \mathfrak{g}$$

and

$$H \setminus X \cong \text{Ad}(H)\text{-conjugacy classes of minimal parabolic subalgebras of } \mathfrak{g}.$$

**Theorem 2 [M-O].** *Let  $G$  and  $H$  be as in the following list (complex classical cases). Then we can express  $H \setminus X$  (and  $H^a \setminus X$ ) by symbols. ( $p+q = n$ ,  $[H : G_0^\sigma] = 1$  or 2.)*

Type	$G$	$H$	$H^a$
AI	$GL(n, \mathbb{C})$	$O(n, \mathbb{C})$	$GL(n, \mathbb{R})$
AII	$GL(n, \mathbb{C})$	$Sp(n/2, \mathbb{C})$ ( $n$ even)	$U^*(n)$
AIII	$GL(n, \mathbb{C})$	$GL(p, \mathbb{C}) \times GL(q, \mathbb{C})$	$U(p, q)$
BI	$SO(2n+1, \mathbb{C})$	$S(O(2p+1, \mathbb{C}) \times O(2q, \mathbb{C}))$	$SO(2p+1, 2q)$
CI	$Sp(n, \mathbb{C})$	$GL(n, \mathbb{C})$	$Sp(n, \mathbb{R})$
CII	$Sp(n, \mathbb{C})$	$Sp(p, \mathbb{C}) \times Sp(q, \mathbb{C})$	$Sp(p, q)$
DI	$SO(2n, \mathbb{C})$	$S(O(2p, \mathbb{C}) \times O(2q, \mathbb{C}))$	$SO(2p, 2q)$
DI'	$SO(2n, \mathbb{C})$	$S(O(2p+1, \mathbb{C}) \times O(2q-1, \mathbb{C}))$	$SO(2p+1, 2q-1)$
DIII	$SO(2n, \mathbb{C})$	$GL(n, \mathbb{C})$	$SO^*(2n)$

*Note.* In [M-O] p.155, we should read  $GL(n, \mathbb{C})$  for  $\mathbb{C}^\times \times PSL(n, \mathbb{C})$  on the line of DIII in Table 1.

Precise description of symbols and many examples are given in [M-O]. But we can explain shortly the essential part as follows.

Let  $x \in U \subset X$ . Then  $\alpha_x = \text{Lie}(P_x) \cap \mathfrak{s}$  is  $\sigma$ -stable by the definition of  $U$ . Let  $\Sigma_x$  be the root system of the pair  $(\mathfrak{g}, \alpha_x)$  and  $\Sigma_x^+$  the positive system of  $\Sigma_x$  corresponding to  $P_x$ . Let  $\Psi_x$  denote the set of simple roots in  $\Sigma_x^+$ . Then we can take an orthogonal basis  $\{e_1, \dots, e_n\}$  of the dual  $\alpha_x^*$  of  $\alpha_x$  such that

$$\Psi_x = \begin{cases} \{\alpha_1, \dots, \alpha_{n-1}\} & \text{if } G = GL(n, \mathbb{C}), \\ \{\alpha_1, \dots, \alpha_n\} & \text{otherwise,} \end{cases}$$

where  $\alpha_1 = e_1 - e_2, \dots, \alpha_{n-1} = e_{n-1} - e_n$  and  $\alpha_n = e_n, 2e_n$  or  $e_{n-1} + e_n$  if  $G = SO(2n+1, \mathbb{C}), Sp(n, \mathbb{C})$  or  $SO(2n, \mathbb{C})$ , respectively.

To the left coset  $(K \cap H)x$  in  $U$ , there corresponds a sequence  $\varepsilon_1 \varepsilon_2 \dots \varepsilon_n$  consisting of the following four kinds of letters.

( $\pm$ ) If  $\sigma e_i = e_i$ , then  $\varepsilon_i = +$  ("a boy") or  $-$  ("a girl"). When  $\varepsilon_i = \pm$  and  $\varepsilon_j = \pm$  ( $i \neq j$ ),

$$\varepsilon_i = \varepsilon_j \iff g(\alpha_x, e_i - e_j) \subset \text{Lie}(H).$$

(a) If  $\sigma e_i = e_j$  with  $i \neq j$ , then we put a small letter ("a family name") to the couple  $(\varepsilon_i, \varepsilon_j)$ .

(A) If  $\sigma e_i = -e_j$  with  $i \neq j$ , then we put a capital letter to the "old" couple  $(\varepsilon_i, \varepsilon_j)$ .

(O) If  $\sigma e_i = -e_i$ , then  $\varepsilon_i = O$  ("the aged" or "dead"?).

Let  $w_i$  be the reflection with respect to the simple root  $\alpha_i$  and  $P_i = P \cup Pw_iP$  ( $P = P_x$ ) the parabolic subgroup of  $G$  for  $\alpha_i$ . Let  $\pi_i$  denote the projection of  $X = G/P$  onto  $G/P_i$ .

**Notation.** For two  $H$ -orbits  $D_1$  and  $D_2$  on  $X$ , we write

$$D_1 \xrightarrow{i} D_2 \iff \pi_i(D_1) = \pi_i(D_2) \text{ and } \dim D_1 < \dim D_2.$$

We put here two examples. (You can see 23 figures of examples in [M-O].)

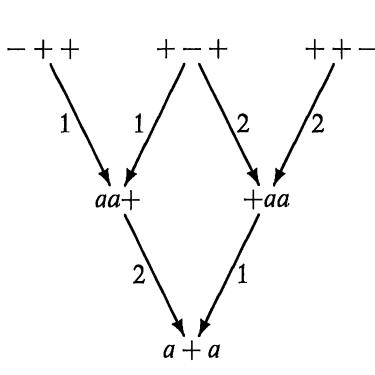


Fig. 1.  $G = GL(3, \mathbb{C})$

$H = GL(2, \mathbb{C}) \times GL(1, \mathbb{C})$

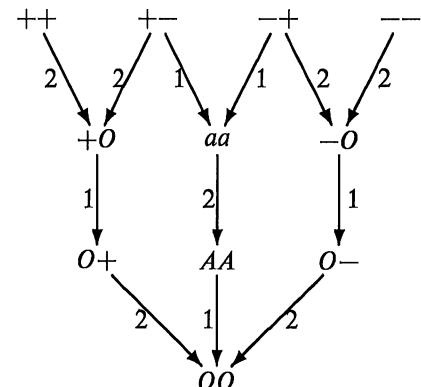


Fig. 2.  $G = Sp(2, \mathbb{C})$ ,  $H = GL(2, \mathbb{C})$

*Remark 2* ([S], [M2]). In complex cases, we can find all the closure relations among  $H$ -orbits on  $X$  from the following two properties.

$$(a) D_1 \xrightarrow{i} D_2 \Rightarrow D_1 \subset D_2^{cl}.$$

$$(b) D_1 \xrightarrow{i} D_2, D_3 \xrightarrow{i} D_4 \text{ and } D_1 \subset D_3^{cl} \Rightarrow D_2 \subset D_4^{cl}.$$

This is proved by the same argument as that of the Bruhat ordering since

$$D_1 \xrightarrow{i} D_2 \text{ and } D_1 \xrightarrow{i} D_3 \Rightarrow D_2 = D_3$$

in complex cases. To find all the closure relations in general real cases, we should follow a rather complicated procedure given in [M2].

*Remark 3.* These diagrams of orbits are useful to the study of the asymptotic behavior of spherical functions on semisimple symmetric spaces ([O]) and embeddings of Harish-Chandra modules into principal series ([M-O]).

*Remark 4 (Problem).* If  $\Sigma = \Sigma(\mathfrak{g}, \mathfrak{a})$  is classical, then there exists (in principle) a similar (sometimes the same) expression of the  $H$ -orbits on  $X$  as that in a complex case. Give a complete list of such expressions by symbols. (For example, it is proved in [M2] that the diagram of  $H^a \setminus X$  is upside-down to that of  $H \setminus X$ .)

*Example 2 (= Exercise).* When  $G = GL(n, \mathbb{F})$  and  $H = GL(p, \mathbb{F}) \times GL(n-p, \mathbb{F})$  for a division algebra  $\mathbb{F}$  of characteristic  $\neq 2$ , the diagram of the  $H$ -orbits on  $X$  does not depend on  $\mathbb{F}$ .

**Problem.** Give good symbols for  $H$ -orbits on  $X$  when  $\Sigma$  is exceptional.

### 3. Uzawa's Function $f$ and Vector Field $v$ on $X$ (Related to Intersections of $H$ -Orbits and $H^a$ -Orbits on $X$ )

Recently, T. Uzawa discovered the following function  $f$  and vector field  $v$  on  $X$  which have very nice properties with respect to  $H$ -orbits and  $H^a$ -orbits.

Let  $Y_0$  be a generic element of  $\mathfrak{s}$ . Then  $Y_0$  defines a minimal parabolic subgroup  $P_0$  of  $G$  such that  $Y_0 \in \mathfrak{a}_0 = \text{Lie}(P_0) \cap \mathfrak{s}$  and that  $Y_0$  is dominant for the positive system of the root system  $\Sigma(\mathfrak{g}, \mathfrak{a}_0)$  corresponding to  $P_0$ . By the natural identification

$$G/P_0 \cong K/M_0 \cong \text{Ad}(K)Y_0$$

( $K \cap P_0 = M_0$  = the centralizer of  $Y_0$  in  $K$ ),  $X = G/P_0$  is embedded into  $\mathfrak{s}$ . Let  $Y_x$  denote the element in  $\text{Ad}(K)Y_0$  corresponding to  $x \in X$ .

**Definition.** (i) We define a function  $f$  on  $X$  by  $f(x) = |Y_x^+|^2 = B(Y_x^+, Y_x^+)$  on  $X$  where  $Y_x^+ = \frac{1}{2}(Y_x + \sigma Y_x)$  and  $B(,)$  is the Killing form on  $\mathfrak{g}$ .

(ii) A vector field  $v$  on  $X$  is defined by  $v_x$  = the (infinitesimal)  $Y_x^+$ -action at  $x$  for  $x \in X$ .

(iii)  $\Phi_t$  ( $t \in \mathbb{R}$ ) is the one-parameter group of transformations of  $X$  for the vector field  $v$ .

(iv)  $\Phi_{\pm\infty}(x) = \lim_{t \rightarrow \pm\infty} \Phi_t(x)$  for  $x \in X$ .

*Remark 5.* The vector field  $v$  is the gradient of the function  $f$  with respect to the  $K$ -invariant Riemannian metric on  $X = K/M_0$  induced from the inner product  $(Z, Z') = B([Z, Y_0], Z'_s)$  on  $\mathfrak{k}^{\perp m_0}$  where  $Z'_s$  is the element in  $\mathfrak{s}$  such that  $Z'_s - Z' \in \text{Lie}(P_0)$  and  $\mathfrak{k}^{\perp m_0}$  is the orthogonal complement of  $m_0 = \text{Lie}(M_0)$  in  $\mathfrak{k}$ .

*Remark 6.* If the real rank of  $G$  is larger than one, then  $f$  and  $v$  depend essentially (not constant multiple) on the choice of  $Y_0$ .

*Example 3.* (continued to Example 1) Take

$$Y_0 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \in \mathfrak{s} = \left\{ \begin{pmatrix} z & x+iy \\ x-iy & -z \end{pmatrix} \mid x, y, z \in \mathbb{R} \right\}.$$

Since  $P_0$  is the subgroup of  $G$  consisting of upper triangular matrices,  $eP_0$  corresponds to  $\infty$  in  $P^1(\mathbb{C}) = \mathbb{C} \cup \{\infty\}$  and

$$kP_0 \mapsto \begin{pmatrix} a & b \\ -\bar{b} & \bar{a} \end{pmatrix} \infty = \frac{a}{-\bar{b}} \quad \text{for } k = \begin{pmatrix} a & b \\ -\bar{b} & \bar{a} \end{pmatrix} \in K.$$

On the other hand,

$$\begin{aligned} \begin{pmatrix} a & b \\ -\bar{b} & \bar{a} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} a & b \\ -\bar{b} & \bar{a} \end{pmatrix}^{-1} &= \begin{pmatrix} a & -b \\ -\bar{b} & -\bar{a} \end{pmatrix} \begin{pmatrix} \bar{a} & -b \\ \bar{b} & a \end{pmatrix} \\ &= \begin{pmatrix} a\bar{a} - b\bar{b} & -2ab \\ -2\bar{a}\bar{b} & -a\bar{a} + b\bar{b} \end{pmatrix}. \end{aligned}$$

So  $\text{Ad}(K)Y_0$  is the sphere given by  $x^2 + y^2 + z^2 = 1$  and the function  $f$  is given by  $z^2$ . Two points  $\{\infty\}$ ,  $\{0\}$  and the unit circle in  $P^1(\mathbb{C})$  correspond to  $(0, 0, 1)$ ,  $(0, 0, -1)$  and the circle defined by  $z = 0$ , respectively, in  $\text{Ad}(K)Y_0$ .

**Theorem 3** [U]. (i)  $v$  is tangent to  $H$ -orbits and  $H^a$ -orbits.

(ii)  $(df)_x = 0 \iff v_x = 0 \iff x$  is special.

(iii) Let  $D$  be an  $H$ -orbit on  $X$ . Then there exists  $m = \min_{x \in D} f(x)$  and for  $x \in D$ ,

$$f(x) = m \iff x \text{ is special}.$$

(iv)  $\Phi_{-\infty}(D) = D \cap U$  for  $H$ -orbits  $D$  on  $X$ .

**Corollary 1** [M3]. (a)  $D \cap D^a = (K \cap H)x$  for an  $x \in U$ .

(b) For two  $H$ -orbits  $D$  and  $E$  on  $X$ ,

$$D^{cl} \supset E \iff D \cap E^a \neq \emptyset \iff D^a \subset (E^a)^{cl}.$$

*Proof* ([U]). (a) Let  $x \in D \cap D^a$ . We have only to show that  $x \in U$  by Theorem 1. Let  $m$  be the value of the function  $f$  at the points in  $D \cap U$  ( $= D^a \cap U$ ). Suppose that  $x \notin U$ . Then  $f(x) > m$  by (iii). Since the function for the  $H^a$ -orbit structure is  $|Y_0|^2 - f(x)$ , we have also  $f(x) < m$  by (iii), a contradiction.

(b) Since  $(H^a)^a = H$ , we have only to prove the left  $\iff$ .

The assertion  $D^{cl} \supset E \Rightarrow D \cap E^a \neq \emptyset$  is clear since

$$T_x(E) + T_x(E^a) = T_x(X)$$

for any  $x \in E \cap E^a$  ([M3]).

Suppose that  $D \cap E^a \neq \emptyset$  and let  $x \in D \cap E^a$ . Then

$$\Phi_\infty(x) = \lim_{t \rightarrow \infty} \Phi_t(x) \in D^{cl} \cap E^a \cap U = D^{cl} \cap E \cap U$$

by (i) and (iv). Hence  $D^{cl} \cap E \neq \emptyset$  and therefore  $D^{cl} \supset E$ .  $\square$

**Corollary 2.** *Let  $D$  be an  $H$ -orbit on  $X$  and  $x \in D \cap D^a$ . Then*

$$(i) \quad D \cong (K \cap H) \times_L \Phi_{-\infty}^{-1}(x)$$

where  $L = K \cap H \cap P_x$  and

$$(ii) \quad D \cap E^a \cong (K \cap H) \times_L (\Phi_{-\infty}^{-1}(x) \cap E^a)$$

for any  $H^a$ -orbit  $E^a$  on  $X$ . (Moreover it is clear that the fibers  $\Phi_{-\infty}^{-1}(x)$  and  $\Phi_{-\infty}^{-1}(x) \cap E^a$  are contractible to the point  $x$ .)

#### 4. Remarks on Spherical Subgroups

Suppose that  $G$  is a complex semisimple Lie group. A complex Lie subgroup  $H$  of  $G$  is called “spherical” if there exists an open  $H$ -orbit on  $X$ . Such pairs  $(G, H)$  are classified by [K] when  $G$  is simple and  $H$  is reductive, and by [Br2] in general.

**Theorem 4** [Br1, V].  *$H \subset G$  is spherical  $\iff \#(H \setminus X)$  is finite. (Note that  $\Leftarrow$  is clear.)*

There is a simple proof of  $\Rightarrow$  using “rank-one sections” as follows.

*Proof.* We may assume that  $HP$  is open in  $G$ . Write  $G = P_{\beta_1}P_{\beta_2} \cdots P_{\beta_m}$  where the  $\beta_i$ 's are simple roots and  $P_{\beta_i} = P \cup Pw_{\beta_i}P$ . Put  $P^{(i)} = P_{\beta_1}P_{\beta_2} \cdots P_{\beta_i}$  ( $P^{(0)} = P$ ). We will show

$$\#(H \setminus HP^{(i)}/P) < \infty \text{ for } i = 0, 1, \dots, m$$

by induction on  $i$ .

By the hypothesis of induction, we may assume that

$$HP^{(i-1)} = Hg_1P \cup \cdots \cup Hg_kP .$$

Then we have

$$HP^{(i)} = Hg_1P_{\beta_i} \cup \cdots \cup Hg_kP_{\beta_i} .$$

We have only to show that  $\#(H \setminus Hg_jP_{\beta_i}/P) < \infty$  for  $j = 1, \dots, k$ . Since  $HP^{(i-1)}$  is open in  $G$ ,  $(g_jP_{\beta_i}/P) \cap (HP^{(i-1)}/P)$  is (Zariski) open in the one-dimensional subvariety  $g_jP_{\beta_i}/P$  of the complex algebraic variety  $X$ . Hence the compliment of  $(g_jP_{\beta_i}/P) \cap (HP^{(i-1)}/P)$  in  $g_jP_{\beta_i}/P$  consists of finite points and therefore  $\#(H \setminus Hg_jP_{\beta_i}/P) < \infty$ .  $\square$

Let  $G$  be a real semisimple Lie group and  $H$  a Lie subgroup of  $G$ .

**Conjecture 1.** *If the real rank of  $G$  is one and there exists an open  $H$ -orbit on  $X = G/P$ , then  $\#(H \setminus X) < \infty$ .*

By the same argument as above for spherical subgroups, Conjecture 1 implies the following Conjecture 2.

**Conjecture 2.** *If there exists an open  $H$ -orbit on  $X$ , then  $\#(H \setminus X) < \infty$ .*

*Remark 7.* In general,  $\#(H \setminus G/P) < \infty$  does not imply  $\#(H_{\mathbb{C}} \setminus G_{\mathbb{C}}/P_{\mathbb{C}}) < \infty$ . For example, if  $G = SU(n, 1)$  ( $n > 2$ ) and  $H = \theta N$  (where  $N$  is the unipotent radical of  $P$ ), then  $\#(H \setminus G/P) = 2$  and  $\#(H_{\mathbb{C}} \setminus G_{\mathbb{C}}/P_{\mathbb{C}}) = \infty$ .

## References

- [A] Aomoto, K.: On some double coset decompositions of complex semi-simple Lie groups. *J. Math. Soc. Japan* **18** (1966) 1–44
- [Be] Berger, M.: Les espaces symmétiques non compacts. *Ann. Sci. École Norm. Sup.* **74** (1957) 85–177
- [Br1] Brion, M.: Quelques propriétés des espaces homogènes sphériques. *Manuscripta Math.* **55** (1986) 191–198
- [Br2] Brion, M.: Classification des espaces homogènes sphériques. *Comp. Math.* **63** (1987) 189–208
- [K] Krämer, M.: Sphärische Untergruppen in Kompakten zusammenhängenden Liegruppen. *Comp. Math.* **38** (1979) 129–153
- [M1] Matsuki, T.: The orbits of affine symmetric spaces under the action of minimal parabolic subgroups. *J. Math. Soc. Japan* **31** (1979) 331–357
- [M2] Matsuki, T.: Closure relations for orbits on affine symmetric spaces under the action of minimal parabolic subgroups. *Adv. Studies Pure Math.* **14** (1988) 541–559
- [M3] Matsuki, T.: Closure relations for orbits on affine symmetric spaces under the action of parabolic subgroups. *Intersections of associated orbits, Hiroshima Math. J.* **18** (1988) 59–67
- [M-O] Matsuki, T., Oshima, T.: Embeddings of discrete series into principal series. In: *The Orbit Method in Representation Theory*. Birkhäuser, Boston 1990, pp. 147–175
- [O] Oshima, T.: Asymptotic behavior of spherical functions on semisimple symmetric spaces. *Adv. Studies Pure Math.* **14** (1988) 561–601
- [R] Rossmann, W.: The structure of semisimple symmetric spaces. *Canad. J. Math.* **31** (1979) 157–180
- [S] Springer, T. A.: Some results on algebraic groups with involutions. *Adv. Studies Pure Math.* **6** (1984) 525–534
- [U] Uzawa, T.: Invariant hyperfunction sections of line bundles. Preprint (1990)
- [V] Vinberg, E. B.: Complexity of actions of reductive groups. *Funct. Anal. Appl.* **20** (1985) 1–11



# Sur les Formes Automorphes de Carré Intégrable

Colette Maeglin

URA 748, UFR de Mathématiques, Université de Paris VII, 2, place de Jussieu  
F-75251 Paris Cedex 05, France

## 1. Définitions et notations

Soit  $k$  un corps global, pour simplifier de caractéristique 0, i.e. une extension algébrique finie de  $\mathbb{Q}$ . On note  $\mathbb{A}$  l'anneau des adèles de  $k$  c'est-à-dire la  $k$ -algèbre formée des éléments  $\prod_v x_v$  où  $v$  parcourt l'ensemble des places de  $k$  et où  $x_v$  est un élément du complété  $k_v$  de  $k$  en la place  $v$ , entier pour tout  $v$  sauf un nombre fini.

Soit  $\mathbf{G}$  un groupe algébrique affine connexe, défini sur  $k$  que l'on suppose réductif. Les exemples classiques s'obtiennent en considérant un espace vectoriel de dimension finie  $V$  défini sur  $k$  et une forme bilinéaire  $\Phi$  sur  $V$  (définie sur  $k$ ); on prend alors pour  $\mathbf{G}$  le groupe des automorphismes linéaires de  $V$  qui respectent  $\Phi$ . Pour toute  $k$ -algèbre,  $k'$ , on note  $\mathbf{G}(k')$  les points de  $\mathbf{G}$  définis sur  $k'$ .

On appelle forme automorphe sur  $\mathbf{G}$ , une fonction  $\phi$  à valeurs complexes sur  $\mathbf{G}(k) \backslash \mathbf{G}(\mathbb{A})$  qui vérifie un certain nombre de propriétés pour lesquelles je renvoie à [3]. On appelle forme automorphe de carré intégrable modulo le centre, une forme automorphe,  $\phi$ , pour laquelle il existe un caractère unitaire  $\omega$  du centre  $\mathbf{Z}(\mathbb{A})$  de  $\mathbf{G}(\mathbb{A})$  tel que:

$$(i) \quad \phi(zg) = \omega(z)\phi(g), \quad \forall z \in \mathbf{Z}(\mathbb{A}), \quad \forall g \in \mathbf{G}(\mathbb{A}),$$

$$(ii) \quad \int_{\mathbf{G}(k)\mathbf{Z}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A})} \phi(g)\overline{\phi}(g) dg < \infty,$$

où  $dg$  est une mesure de Haar sur  $\mathbf{G}(k)\mathbf{Z}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A})$ .

Dans cet exposé, on s'intéresse à la détermination des formes automorphes de carré intégrable. L'article de base pour cette question est [5] (cf. aussi [8]).

On appelle sous-groupe parabolique de  $\mathbf{G}$ , un sous-groupe fermé  $\mathbf{P}$  de  $\mathbf{G}$  tel que  $\mathbf{G}(\bar{k})/\mathbf{P}(\bar{k})$  soit une variété projective ( $\bar{k}$  est une clôture algébrique de  $k$ ). Ces groupes sont (comme  $\mathbf{G}$ ) connexes et jouent un rôle privilégié dans la théorie. Si  $\mathbf{P} \neq \mathbf{G}$ ,  $\mathbf{P}$  n'est plus un groupe réductif mais comme tout groupe algébrique affine connexe, il contient un unique sous-groupe fermé normal, noté " $\mathbf{P}$ " tel que  $\mathbf{P}/"\mathbf{P}$  soit réductif et " $\mathbf{P}$ " est minimal avec cette propriété. On pose

$$\mathbf{M}_P := \mathbf{P}/"\mathbf{P}.$$

On note  $X_P$  le groupe des caractères de  $\mathbf{M}_P(\mathbb{A})$ , à valeurs complexes, engendré par les caractères du type  $|\chi|^s$  où  $||$  est la valeur absolue adélique, où  $\chi$  est

un caractère rationnel de  $\mathbf{M}_P$  et où  $s$  est un nombre complexe. Tout caractère appartenant à  $X_P$  est évidemment trivial sur  $\mathbf{M}_P(k)$ . On vérifie que  $X_P$  est un espace vectoriel.

Par exemple si  $\mathbf{G} = \text{Aut}(V, \Phi)$  avec  $\Phi \equiv 0$ , i.e.  $\mathbf{G} \simeq \text{GL}(V)$  alors les sous-groupes paraboliques de  $\mathbf{G}$  sont les stabilisateurs des drapeaux:

$$V_0 = 0 \subset V_1 \subset V_2 \subset \cdots \subset V_d = V$$

de  $V$ . Ici  $\mathbf{M}_P \simeq \prod_{1 \leq i \leq d} \text{GL}(V_i/V_{i-1})$  et  $X_P \simeq \mathbb{C}^d$  isomorphisme donné par l'application qui à  $(z_1, \dots, z_d) \in \mathbb{C}^d$  associe le caractère  $\chi((m_1, \dots, m_d)) = |\det m_1|^{z_1} \cdots |\det m_d|^{z_d}$ .

Soit  $\mathbf{P}$  un sous-groupe parabolique de  $\mathbf{G}$ ; on généralise la notion de formes automorphes en notant  $A(\mathbf{P}(k) \backslash \mathbf{P}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A}))$  l'ensemble des fonctions à valeurs complexes sur ce quotient qui vérifient les mêmes propriétés que les formes automorphes. En fait pour tout  $\phi \in A(\mathbf{P}(k) \backslash \mathbf{P}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A}))$  et pour tout  $g \in \mathbf{G}(\mathbb{A})$  la fonction  $\phi_g$  sur  $\mathbf{M}_P(\mathbb{A})$  définie par:

$$\phi_g(m) = \phi(mg)\delta_P(m)^{-1/2}$$

est une forme automorphe pour  $\mathbf{M}_P$  ( $\delta_P(m)$  est le Jacobien de l'automorphisme de  $\mathbf{P}(\mathbb{A})$  qui envoie  $u \in \mathbf{P}(\mathbb{A})$  sur  $um^{-1}$ ).

On aura aussi besoin de la notion de formes automorphes cuspidales; soit  $\phi$  comme ci-dessus, on demande en plus que l'on ait pour tout sous-groupe parabolique propre  $\mathbf{P}'$  de  $\mathbf{M}_P$ :

$$\forall g \in \mathbf{G}(\mathbb{A}), \quad \int_{\mathbf{P}'(k) \backslash \mathbf{P}'(\mathbb{A})} \phi(ug) du = 0.$$

Cela revient à dire que la représentation de  $\mathbf{M}_P(\mathbb{A})$  engendrée par  $\phi_g$  (aux places à l'infini, c'est un  $\mathbf{g} - K$  module) est une représentation automorphe cuspidale.

## 2. Séries d'Eisenstein

Les séries d'Eisenstein fournissent un procédé pour transformer un élément  $\phi \in A(\mathbf{P}(k) \backslash \mathbf{P}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A}))$  en une famille,  $E(\phi, \lambda)$  de formes automorphes dépendant méromorphiquement de  $\lambda \in X_P$ . Toutefois l'existence de  $E(\phi, \lambda)$  n'a de démonstration écrite que si  $\phi$  est de carré intégrable modulo le centre (cf. [5], Chapitre 6 et 7) (Le résultat général est annoncé par Bernstein et dans le cas, que nous ne considérons pas d'un corps global de caractéristique  $> 0$ , résulte de résultats de [11] et [15]).

Dans ce qui suit, on suppose que  $\phi \in A(\mathbf{P}(k) \backslash \mathbf{P}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A}))$  vérifie la condition (i) ci-dessus (i.e. à un caractère central unitaire) pour l'action du centre de  $\mathbf{M}_P(\mathbb{A})$  et est cuspidale. Alors  $\phi$  est de carré intégrable modulo le centre de  $\mathbf{M}_P(\mathbb{A})$ . On sait que les pôles de  $E(\phi, \lambda)$  sont de nature très simple: soit  $\lambda_0 \in X_P$ , alors il existe un ensemble fini  $E$  d'hyperplans de  $X_P$  passant par  $\lambda_0$  et pour chacun de ces hyperplans  $H \in E$  un entier  $n_H$  ( $E$  et  $n_H$  peuvent être choisis indépendamment de  $\phi$ ) tels que:

$$\prod_{H \in E} P_H^{n_H} E(\phi, \lambda)$$

soit holomorphe en  $\lambda_0$  où  $P_H$  est l'équation de l'hyperplan  $H$ . On peut définir les résidus successifs de  $E(\phi, \lambda)$  le long d'un ensemble fini ordonné,  $E'$ , d'hyperplans linéairement indépendants. On obtient une famille de formes automorphes sur  $\mathbf{G}(\mathbb{A})$  dépendant méromorphiquement de  $\lambda$  appartenant à  $\cap_{H \in E'} H$ .

### Nature des pôles des Séries d'Eisenstein

Admettons que l'on sache normaliser les opérateurs d'entrelacements (cf. [5], appendice 2 et [14]); on sait en particulier le faire si  $\mathbf{G} \simeq \mathrm{GL}(V)$ . Les pôles des séries d'Eisenstein sont alors de 3 types différents :

- le premier type est de nature locale, il correspond aux pôles des opérateurs d'entrelacements normalisés locaux;
- les deux autres types sont de nature globale : les facteurs de normalisations des opérateurs d'entrelacements sont des quotients de fonctions  $L$ . Le deuxième type correspond aux zéros des fonctions  $L$  qui apparaissent aux dénominateurs et le troisième type correspond aux pôles des fonctions  $L$  qui apparaissent aux numérateurs.

### Forme faible des résultats de Langlands

**Théorème (Langlands).** *Toute forme automorphe de carré intégrable modulo le centre est un résidu de séries d'Eisenstein.*

En fait la théorie de Langlands [5] donne beaucoup plus de renseignements puisque elle limite très sérieusement les hyperplans qui peuvent fournir des résidus de carré intégrable.

**Conjecture.** *Soient  $\phi, \lambda, E(\phi, \lambda)$  comme ci-dessus. Les seuls hyperplans singuliers pour  $E(\phi, \lambda)$  qui peuvent fournir des résidus de carré intégrable sont les hyperplans singuliers provenant des pôles des fonctions  $L$  des numérateurs des facteurs de normalisations.*

## 3. Le cas de $\mathrm{GL}(V)$

La conjecture est vraie et on démontre alors aisément le théorème suivant (cf. [9], conjecturé et démontré dans des cas particulier par Jacquet [4] (la partie (i) avait été obtenue par Speh [13])) :

**Théorème.** *Soit  $V$  un espace vectoriel de dimension  $n$  sur  $k$ ; on suppose que  $\mathbf{G} = \mathrm{GL}(V)$ .*

(i) *Soit  $n = da$  une décomposition de  $n$  et soit*

$$V_0 = 0 \subset \cdots \subset V_i \subset \cdots \subset V_a = V$$

*un drapeau tel que  $\dim V_i = id$ , pour  $0 \leq i \leq a$ . On note  $\mathbf{P}$  le stabilisateur de ce drapeau et alors  $\mathbf{M}_P \simeq \prod_{1 \leq i \leq a} \mathrm{GL}(d_i)$ . Soit  $\phi \in A(\mathbf{P}(k))^\ast \mathbf{P}(\mathbb{A}) \backslash \mathbf{G}(\mathbb{A})$  telle que pour tout  $g \in \mathbf{G}(\mathbb{A})$  la représentation automorphe de  $\mathbf{M}_P(\mathbb{A})$  engendrée par  $\phi_g$  soit*

*irréductible isomorphe à  $\varrho \otimes \cdots \otimes \varrho$  où  $\varrho$  est une représentation automorphe cuspidale fixée de  $\mathrm{GL}(d)(\mathbb{A})$  de caractère central unitaire. On rappelle que  $X_P \simeq \mathbb{C}^a$ . Alors l'ensemble des hyperplans siuguliers pour  $E(\phi, \lambda)$  passant par le point:*

$$(a-1)/2, \dots, (a-2i+1)/2, \dots, -(a-1)/2$$

*est l'ensemble des hyperplans  $s_i - s_{i-1} = 1$  pour  $1 \leq i \leq a$ . Le résidu obtenu successivement le long de ces hyperplans rangés dans n'importe quel ordre (l'ordre est indifférent) est une forme automorphe sur  $\mathbf{G}(\mathbb{A})$  de carré intégrable modulo le centre.*

(ii) *Toutes les formes automorphes de carré intégrable modulo le centre sont obtenues de cette façon.*

*Sur la preuve.* (i) C'est un problème combinatoire peu profond (déjà résolu dans l'appendice 3 de [5] pour  $n = 4$ ) que de montrer que les zéros des fonctions  $L$  des dénominateurs des facteurs de normalisations n'interviennent pas.

(ii) C'est un problème de décomposition de certaines induites (du moins il faut calculer leurs quotients irréductibles) qui permet de régler le cas des pôles des opérateurs d'entrelacements locaux.

## 4. Interprétation d'Arthur

Dans [6] Langlands suggère que les classes d'isomorphie de représentations automorphes tempérées se regroupent en paquets qui devraient être paramétrés par les classes de conjugaison d'homomorphismes "admissibles" d'un groupe (dont l'existence devrait provenir de la théorie des catégories tanakiennes), noté  $L_F$ , dans le  $L$ -groupe associé à  $G$ . Admettant cela, Arthur suggère que les classes d'isomorphie de représentations automorphes de carré intégrable se regroupent en paquets que devraient être paramétrés par les classes de conjugaison d'homomorphismes "admissibles" du produit direct de  $L_F$  par  $\mathrm{SL}(2, \mathbb{C})$  dans le  $L$ -groupe. Arthur donne aussi une paramétrisation (conjecturale) des représentations intervenant dans un paquet donné ainsi que de la multiplicité avec laquelle elles devraient intervenir.

En admettant la conjecture de Ramanujan, i.e. que l'ensemble des formes automorphes cuspidales de  $\mathrm{GL}(n)$ , de caractère central unitaire, est exactement l'ensemble des formes automorphes tempérées, le résultat cité pour  $\mathrm{GL}(n)$  prouve les conjectures d'Arthur dans ce cas (chaque paquet est ici réduit à un élément et la multiplicité est un).

Un cas particulier intéressant des conjectures d'Arthur est de considérer les homomorphismes "admissibles" triviaux sur  $L_F$ . Les paquets qu'ils devraient paramétriser sont appelés, par Arthur, unipotents. Supposons pour simplifier grandement que  $G$  est un groupe classique déployé. L'ensemble des paquets unipotents devrait alors être en bijection avec l'ensemble des orbites unipotentes du groupe complexe dual,  $G^*$ , associé à  $G$  (i.e. la composante neutre du  $L$ -groupe) ne rencontrant aucun sous-groupe de Levi propre de  $G^*$ . La méthode employée pour  $\mathrm{GL}(n)$  devrait conduire aisément au résultat suivant:

*soit  $G$  un groupe classique déployé de centre fini; alors les résidus de carré intégrable des séries d'Eisenstein construites à partir des caractères non ramifiés d'un sous-groupe de Borel de  $G$  doivent réaliser avec multiplicité un les*

*représentations automorphes de carré intégrable de  $G(\mathbb{A})$  paramétrées par les orbites unipotentes,  $O$ , de  $G^*$  ne rencontrant aucun sous-groupe de Levi propre de  $G^*$  et par les produits,*

$$\prod_v \chi_v,$$

*indéxés par l'ensemble des places de  $k$ , de caractères du groupe,  $\overline{A}(O)$ , introduit par Lusztig ([7], chapitre 14, où il est noté  $\overline{A}(u)$ ), vérifiant :*

*pour presque tout  $v$ ,  $\chi_v$  est le caractère identité et le caractère de  $\overline{A}(O)$  que l'on obtient en faisant le produit de tous les  $\chi_v$  est le caractère identité.*

Pour un résultat dans ce sens, confer [10]. Ce qui me manque pour obtenir toutes les représentations d'un paquet unipotent est une construction générale des représentations automorphes cuspidales unipotentes; les séries theta en fournit dans certains cas (cf. [12]).

## Bibliographie

1. Arthur, J.: Unipotent automorphic representations: conjectures, in Orbites unipotentes et représentations II. Groupes p-adiques et réels. Astérisque **171-172** (1989) 13–72
2. Arthur, J.: Unipotent automorphic representations: global motivations. In: automorphic forms, Shimura varieties, and L-functions, vol. 1 (L. Clozel et J.S. Milne, eds.). Perspectives in Mathematics, 1990, pp. 1–75
3. Borel, A., Jacquet, H.: Automorphic forms and automorphic representations. In: Proc. of Symp. in Pure Math., vol. 33, vol. 1. Amer. Math. Soc., Providence, RI, 1979, pp. 189–202
4. Jacquet, H.: On the residual spectrum of  $GL(n)$ . In: Lie Group Representations II. (Lecture Notes in Mathematics, vol. 1041 (édité par R. Herb, S. Kudla, R. Lipsmann and J. Rosenberg)) Springer, Berlin Heidelberg New York 1984, pp. 185–208
5. Langlands, R.P.: On the functional equations satisfied by Eisenstein series. (Lecture Notes in Mathematics, vol. 544.) Springer, Berlin Heidelberg New York 1976
6. Langlands, R.P.: Automorphic representations, Shimura varieties and motives – ein Märchen. In: Proc. of Symp. in Pure Math., vol. 33, vol. 2. Amer. Math. Soc., Providence, RI, 1979, pp. 205–246
7. Lusztig, G.: Characters of reductive groups over a finite field. Ann. Math. Studies, vol. 107. Princeton Univ. Press, 1984
8. Mœglin, C., Waldspurger, J.-L.: Décomposition spectrale et séries d'Eisenstein, paraphrase sur l'Écriture. Prépublication, Paris 1989
9. Mœglin, C., Waldspurger, J.-L.: Le spectre résiduel de  $GL(n)$ . Ann. de l'ENS, 1989
10. Mœglin, C.: Orbites unipotentes et spectre discret non ramifié, à paraître à Compositio
11. Morris, L.E.: Eisenstein series for reductive groups over global function fields 1 et 2. Can. J. Math. **34** (1982) 1112–1182 et Can. J. Math. **35** (1983) 974–985
12. Piatetski-Shapiro, I.: On the Saito Kurokawa lifting. Math. inv. **71** (1983) 309–338
13. Speh, B.: Some results on principal series for  $GL(n, \mathbb{R})$ . Ph. D. dissertation M.I.T., Juin 1977
14. Shahidi, F.: A proof of Langlands conjecture on Plancherel measures; complementary series for  $p$ -adic groups. Ann. Math., à paraître
15. Waldspurger, J.-L.: Formes automorphes et séries d'Eisenstein sur un corps de fonctions. Prépublication, Paris 1989



# Semi-simple Groups and Arithmetic Subgroups

*Gopal Prasad*

Tata Institute of Fundamental Research, Homi Bhabha Road  
Colaba, Bombay 400 005, India

*To my teacher  
M.S. Raghu Nath*

In this lecture I shall report on the recent results, and open questions, related to the congruence subgroup problem, computation of the covolume of  $S$ -arithmetic subgroups, bounds for the class-number of simply connected semi-simple groups and state the finiteness theorems of [3]. We shall also briefly mention the recent work on super-rigidity of *cocompact* discrete subgroups of  $\mathrm{Sp}(n, 1)$  and the R-rank 1 form of type  $F_4$ , which implies arithmeticity of these discrete subgroups.

*Notation.* Throughout this report  $k$  is a global field, that is either a number field (i.e. a finite extension of the field  $\mathbf{Q}$  of rational numbers) or the function field of an algebraic curve over a finite field. Let  $V$  be the set of places of  $k$ ,  $V_\infty$  (resp.  $V_f$ ) be the set of archimedean (resp. nonarchimedean) places. For  $v \in V$ ,  $k_v$  will denote the completion of  $k$  at  $v$  with the natural locally compact topology and  $|\cdot|_v$  the normalized absolute value on  $k_v$ . For  $v \in V_f$ ,  $\mathfrak{o}_v$  will denote the ring of integers of  $k_v$ ,  $\mathfrak{f}_v$  the residue field,  $p_v$  the characteristic of  $\mathfrak{f}_v$  and  $q_v$  its cardinality. In the sequel  $k_v$  is assumed to carry the “normalized” Haar measure, see [26, 0.1]. For any finite set  $S$  of places of  $k$  containing  $V_\infty$ ,  $\mathfrak{o}_S$  will denote the ring of  $S$ -integers of  $k$ , i.e.

$$\mathfrak{o}_S = \{x \in k \mid |x|_v \leq 1 \text{ for all } v \notin S\}.$$

$A$  will denote the ring of adèles of  $k$ . For a finite set  $S$  of places of  $k$ , let  $A_S$  be the ring of  $S$ -adèles i.e. the restricted direct product of the  $k_v$ 's for  $v \notin S$ .

Let  $G$  be a connected semi-simple algebraic group defined over  $k$ . We fix an embedding of  $G$  in  $\mathrm{SL}_n$  defined over  $k$  and view  $G$  as a  $k$ -subgroup of  $\mathrm{SL}_n$  in terms of this embedding. Let  $S$  be a fixed finite set of places of  $k$  containing  $V_\infty$ . Let  $G_S = \prod_{v \in S} G(k_v)$  with the locally compact topology induced by the topologies on  $k_v, v \in S$ . We shall let  $\Gamma$  denote the group  $G(k) \cap \mathrm{SL}_n(\mathfrak{o}_S)$ . Note that  $\Gamma$  depends on the embedding of  $G$  in  $\mathrm{SL}_n$  fixed above. Embedded diagonally in  $G_S$ ,  $\Gamma$  is a discrete subgroup of finite covolume. A subgroup of  $G_S$  is said to be  *$S$ -arithmetic* if it is commensurable with  $\Gamma$ .

## 1. The Congruence Subgroup Problem

For any non-zero ideal  $\alpha$  of  $\mathfrak{o}_S$ , we have the “reduction mod  $\alpha$ ”

$$\pi_\alpha : \mathrm{SL}_n(\mathfrak{o}_S) \rightarrow \mathrm{SL}_n(\mathfrak{o}_S/\alpha).$$

The kernel of  $\pi_\alpha|_{\Gamma}$  will be denoted by  $\Gamma_\alpha$ , it is by definition the *principal S-congruence subgroup of  $\Gamma$  of level  $\alpha$* . Since  $\mathfrak{o}_S/\alpha$  is finite,  $\mathrm{SL}_n(\mathfrak{o}_S/\alpha)$  is finite and hence  $\Gamma_\alpha$  is of finite index in  $\Gamma$ . An  $S$ -arithmetic subgroup is an  $S$ -congruence subgroup if it contains a principal  $S$ -congruence subgroup of  $\Gamma$  of some level. It is not difficult to see that this notion (of  $S$ -congruence subgroups) does not depend on the choice of the  $k$ -embedding of  $G$  in  $\mathrm{SL}_n$ .

Henceforth,  $G$  will be assumed to be absolutely almost simple and simply connected. We shall assume further that  $G_S$  is noncompact or, equivalently, for some  $v$  in  $S$ ,  $G$  is isotropic at  $v$  ([24]).

The congruence subgroup problem in its simplest form asks whether any  $S$ -arithmetic subgroup is an  $S$ -congruence subgroup. If the answer is in the affirmative, we say that  $G$  has the congruence subgroup property (for  $S$ -arithmetic subgroups). In general the answer to the above question is in the negative. For example, as has been known since 1880, the group  $\mathrm{SL}_2/\mathbb{Q}$  does not have the congruence subgroup property for  $S = V_\infty$  (but the group  $\mathrm{SL}_n/\mathbb{Q}$  has the congruence subgroup property for all  $n > 2$  – this was proved by Bass-Lazard-Serre and Mennicke independently in 1963). If  $k$  is a totally imaginary number field, the group  $\mathrm{SL}_n/k$  fails to have the congruence subgroup property for any  $n$  ( $S = V_\infty$ ); see [2]. To give a precise measure of the failure, J-P. Serre introduced “the  $S$ -congruence kernel” which is a profinite group defined as follows. On  $G(k)$  we introduce the following two translation invariant topologies:

- (1) The  *$S$ -congruence topology*: In this the  $S$ -congruence subgroups constitute a neighborhood base at the identity. It is obvious that this is the same topology as the one induced on  $G(k)$  from  $G(A_S)$ . By strong approximation ([23], [14], [21]), the completion of  $G(k)$  with respect to the  $S$ -congruence topology is  $G(A_S)$ .
- (2) The  *$S$ -arithmetic topology*: In this the  $S$ -arithmetic subgroups contained in  $G(k)$  constitute a neighborhood base at the identity. Completion of  $G(k)$  with respect to this topology will be denoted by  $\widehat{G}_S$ .

As every  $S$ -congruence subgroup is  $S$ -arithmetic, the  $S$ -arithmetic topology on  $G(k)$  is finer than the  $S$ -congruence topology and therefore there is a continuous homomorphism  $\widehat{G}_S \rightarrow G(A_S)$ . It is not difficult to show that this homomorphism is surjective and its kernel, denoted  $C(S, G)$ , is a profinite group.  $C(S, G)$  is by definition the  *$S$ -congruence kernel*. It is clear that  $C(S, G)$  is trivial if, and only if,  $G$  has the congruence subgroup property (for  $S$ -arithmetic subgroups). In the more precise formulation due to Serre, the congruence subgroup problem is the problem of determining the  $S$ -congruence kernel  $C(S, G)$ .

We have the following topological extension

$$(*) \quad 1 \rightarrow C(S, G) \rightarrow \widehat{G}_S \rightarrow G(A_S) \rightarrow 1$$

of  $G(A_S)$  by  $C(S, G)$ . The natural inclusion of  $G(k)$  in  $\widehat{G}_S$  provides a splitting of this extension over  $G(k) \hookrightarrow G(A_S)$ . It has been conjectured that, under a fairly general hypothesis, (\*) is a central extension i.e.,  $C(S, G)$  is central in  $\widehat{G}_S$ , see Section 4 below. We shall devote the next two sections to topological central extensions.

## 2. Topological Fundamental Group

A topological extension

$$(+) \quad 1 \rightarrow \mathcal{C} \rightarrow \mathcal{E} \rightarrow \mathcal{G} \rightarrow 1 ,$$

of a locally compact and second countable topological group  $\mathcal{G}$  by  $\mathcal{C}$ , with  $\mathcal{E}$  locally compact and second countable, is a universal topological central extension (u.t.c.e.) of  $\mathcal{G}$  if it is a central extension i.e.,  $\mathcal{C}$  is a closed central subgroup of  $\mathcal{E}$ , and given any topological central extension

$$1 \rightarrow C \rightarrow E \rightarrow \mathcal{G} \rightarrow 1 ,$$

with  $E$  locally compact and second countable, there is a *unique* continuous homomorphism  $\varphi : \mathcal{E} \rightarrow E$  making the following diagram commutative

$$\begin{array}{ccccccc} 1 & \rightarrow & \mathcal{C} & \rightarrow & \mathcal{E} & \rightarrow & \mathcal{G} & \rightarrow & 1 \\ & & \downarrow & & \downarrow \varphi & & \parallel & & \\ 1 & \rightarrow & C & \rightarrow & E & \rightarrow & \mathcal{G} & \rightarrow & 1 . \end{array}$$

It is clear that if  $\mathcal{G}$  admits a u.t.c.e., the latter is unique upto natural equivalence. In case (+) is a u.t.c.e. of  $\mathcal{G}$ ,  $\mathcal{C}$  is by definition the *topological fundamental group* of  $\mathcal{G}$  and it is denoted by  $\pi_1(\mathcal{G})$ . If  $\mathcal{G}$  is a connected real semi-simple Lie group, then  $\pi_1(\mathcal{G})$  coincides with the usual (algebraic topological) fundamental group of  $\mathcal{G}$ . It follows from certain results of Moore [20], that if  $\mathcal{G}$  is *perfect* i.e. if it is its own commutator, and the cohomology group  $H_m^2(\mathcal{G}, \mathbf{R}/\mathbf{Z})$ , based on measurable cochains, is finite, then  $\mathcal{G}$  admits a u.t.c.e. and  $\pi_1(\mathcal{G})$  is isomorphic to the dual of  $H_m^2(\mathcal{G}, \mathbf{R}/\mathbf{Z})$ . It is also known that if  $\mathcal{G}$  is totally disconnected, then the cohomology theory of  $\mathcal{G}$  based on measurable cochains is identical with the theory based on continuous cochains [37].

If  $v$  is a nonarchimedean place where  $G$  is isotropic, then  $G(k_v)$  is perfect (in fact any proper normal subgroup is central) and the cohomology group  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$ , defined in terms of continuous cochains, is essentially known:

**Theorem 1.** *Let  $v$  be a nonarchimedean place of  $k$  such that  $G$  is isotropic at  $v$  (or, equivalently,  $G(k_v)$  is noncompact), then  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$  is isomorphic to a subgroup, of index at most two, of the dual  $\hat{\mu}(k_v)$  of the finite group  $\mu(k_v)$  of roots of unity in  $k_v$ . Moreover, if at least one of the following three conditions holds, then it is isomorphic to  $\hat{\mu}(k_v)$ .*

- (i)  $G$  is quasi-split over an odd degree extension of  $k_v$ ;
- (ii)  $k_v$  is not an extension of  $\mathbf{Q}_2$ ;
- (iii)  $k_v$  contains a primitive fourth root of unity.

As a consequence, if  $G$  is isotropic at  $v$ , then  $G(k_v)$  admits a u.t.c.e. and  $\pi_1(G(k_v))$  is isomorphic to the dual of  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$ .

It is expected that for any nonarchimedean place  $v$  where  $G$  is isotropic,  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$  is isomorphic to  $\hat{\mu}(k_v)$ . For the spin group of a quadratic form over  $k$  which is of Witt index at least 2 at  $v$ , this is proved in [27, 1.9] and the same proof would take care of some other classical groups.

For the group  $\mathrm{SL}_2$  the above theorem is due to Moore [19]. For other Chevalley groups (i.e. groups which split over  $k$ ) he proved that  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$  is isomorphic to a subgroup of  $\hat{\mu}(k_v)$ , and about ten years later Deodhar [7] showed that this holds also when  $G$  is quasi-split over  $k_v$  (i.e. contains a Borel subgroup defined over  $k_v$ ). Soon after Moore proved his result, Matsumoto showed, by constructing a suitable topological central extension of  $G(k_v)$ , that if  $G$  is a Chevalley group, the above cohomology group is actually isomorphic to  $\hat{\mu}(k_v)$ , and an observation of Deligne implies that this is also the case if  $G$  is quasi-split over  $k_v$ , see [28, §5]. Bak and Rehmann [1] have proved the above theorem, as well as Theorem 3 stated below, for groups of inner type A of relative rank  $\geq 2$  using  $K$ -theoretic methods.

In the generality stated above, the theorem is proved in [28] using the results of Moore, Matsumoto; Deodhar and Deligne and the Bruhat-Tits theory of reductive groups over nonarchimedean local fields. The complete proof of the above theorem is quite long and difficult and involves some case considerations. It is desirable to have a shorter and simpler proof.

The known results on  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$  in case  $G(k_v)$  is compact (or, equivalently,  $G$  is anisotropic at  $v$ ) are summarised below.

**Theorem 2.** *Let  $v$  be a nonarchimedean place such that  $G(k_v)$  is compact (then, as is well known, there is a central division algebra  $D_v$  over  $k_v$  such that  $G(k_v)$  is isomorphic to the group  $\mathrm{SL}_1(D_v)$  of elements of reduced norm 1 in  $D_v$ , and) the cohomology group  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$ , based on continuous cochains, is a finite group of order a power of  $p_v$ , where  $p_v$  is the characteristic of the residue field of  $k_v$ . It is cyclic if  $D_v$  is not the quaternion central division algebra over  $\mathbf{Q}_2$  and is trivial if  $k_v$  does not contain a primitive  $p_v$ -th root of unity and  $D_v$  is not the quaternion central division algebra over  $\mathbf{Q}_3$ .*

This theorem is proved in [30]. The precise computation of  $H^2(\mathrm{SL}_1(D_v), \mathbf{R}/\mathbf{Z})$  has not yet been done. We conjecture that it is isomorphic to  $\mathbf{Z}/2\mathbf{Z} \oplus \mathbf{Z}/2\mathbf{Z}$  if  $D_v$  is the quaternion central division algebra over  $\mathbf{Q}_2$ , it is isomorphic to  $\mathbf{Z}/3\mathbf{Z}$  if  $D_v$  is the quaternion central division algebra over  $\mathbf{Q}_3$ , and is isomorphic to the  $p_v$ -primary component of the dual  $\hat{\mu}(k_v)$  of the group of roots of unity in  $k_v$  in all other cases.

*Remark.* Theorems 1 and 2 imply that, for any finite set  $S$  of places of  $k$ ,  $H^2(G(A_S), \mathbf{R}/\mathbf{Z})$  is the direct product of the  $H^2(G(k_v), \mathbf{R}/\mathbf{Z})$ ,  $v \notin S$ . If moreover  $G(A_S)$  is perfect, then it admits a u.t.c.e. and its topological fundamental group is the direct sum (with discrete topology) of the  $\pi_1(G(k_v))$ ,  $v \notin S$ ; see [27, Theorem 2.4]. This implies that if  $C(S, G)$  is central in  $\widehat{G}_S$ , then it is actually finite [27, §2].

### 3. The $S$ -Metaplectic Kernel

Let  $\mathcal{H}$  be a subgroup of a locally compact second countable topological group  $\mathcal{G}$ . Assume that  $\mathcal{H}$  is perfect and  $\mathcal{G}$  admits a u.t.c.e. Then there is a topological central extension

$$1 \rightarrow C \rightarrow E \rightarrow \mathcal{G} \rightarrow 1,$$

with  $E$  locally compact and second countable, which splits over  $\mathcal{H}$  and which is universal with respect to this property. The *relative topological fundamental group*  $\pi_1(\mathcal{G}, \mathcal{H})$  is then by definition the group  $C$ .

The  $S$ -metaplectic kernel is the group

$$M(S, G) = \text{Ker}\left(H^2(G(A_S), \mathbf{R}/\mathbf{Z}) \xrightarrow{\text{rest}} H^2(G(k), \mathbf{R}/\mathbf{Z})\right);$$

where  $H^2(G(k), \mathbf{R}/\mathbf{Z})$  denotes the second cohomology of the abstract group  $G(k)$  with coefficients  $\mathbf{R}/\mathbf{Z}$ . The topological central extensions of  $G(A_S)$  by  $\mathbf{R}/\mathbf{Z}$ , which split over the subgroup  $G(k)$ , are classified by the  $S$ -metaplectic kernel. It is obvious that if  $G(k)$  is perfect, then  $M(S, G)$  is isomorphic to the Pontrjagin dual of the relative fundamental group  $\pi_1(G(A_S), G(k))$ .

We now come back to the congruence subgroup problem. Assume that  $G(k)$  is perfect and  $C(S, G)$  is central in  $\widehat{G}_S$  (see Section 4 below). Then adapting an argument of [2, §15] and using Theorem 1, it can be proved that  $(*)$  is the universal extension in the category of topological central extensions of  $G(A_S)$  splitting over  $G(k)$ , see [27, §2]. In particular,  $C(S, G)$  is isomorphic to the relative fundamental group  $\pi_1(G(A_S), G(k))$  and so it is isomorphic to the Pontrjagin dual of the  $S$ -metaplectic kernel  $M(S, G)$ . Thus to determine the  $S$ -congruence kernel, it is enough to compute  $M(S, G)$ . Also, in the theory of automorphic forms (of fractional weights) it is of critical importance to know the topological central extensions of  $G(A)$  which split over  $G(k)$ ; these are determined by  $M(\phi, G)$ . Now we state the following theorem which “determines”  $M(S, G)$  for all  $k$ -isotropic  $G$ .

**Theorem 3.** *Assume that  $G$  is isotropic over  $k$ . Let  $S$  be an arbitrary finite set of places of  $k$ . Then  $M(S, G)$  is trivial if  $S$  contains either a nonarchimedean place, or a real place  $v$  such that the group  $G(k_v)$  is not topologically simply connected. In general  $M(S, G)$  is isomorphic to a subgroup of the dual  $\widehat{\mu}(k)$  of the group of roots of unity in  $k$ .*

For Chevalley groups this theorem was proved by Moore [19] and for groups which are quasi-split over  $k$ , it was proved by Deodhar [7]. For the group  $G = \text{SL}_2$ , Moore in fact proved that  $M(S, G)$  is trivial if  $S$  contains a noncomplex place and it is isomorphic to  $\widehat{\mu}(k)$  otherwise<sup>1</sup>. Soon after this, Matsumoto proved that for all Chevalley groups  $G$ ,  $M(S, G)$  is isomorphic to  $\widehat{\mu}(k)$  if  $S$  does not contain any noncomplex place. The same holds for any group which is quasi-split over  $k$  as was observed by Deligne. If either  $S \supset V_\infty$  or  $k$  is a totally imaginary number field, the precise computation of  $M(S, G)$  for the groups  $\text{SL}_n (n \geq 3)$  and

<sup>1</sup> This result is equivalent to his theorem on the “uniqueness of the reciprocity law of global class field theory” – see [4] for an elegant proof of the latter.

$\mathrm{Sp}_{2n}(n \geq 2)$  is already in [2], and following the ideas of this paper, Vaserstein [35] computed the metaplectic kernel for many other classical groups. In 1981, Bak, in a Comptes Rendus note, outlined a proof of this theorem for all groups of classical type of relative rank at least two which uses the results of [1].

For arbitrary simply connected  $k$ -isotropic groups, the above theorem was proved by Prasad and Raghunathan in 1980 [27], and besides the results of Moore and Deodhar for split and quasi-split groups, the proof uses the results of [28] on topological central extensions of  $G(k_v)$ . Note that for a real place  $v$ , the condition that  $G(k_v)$  is not topologically simply connected is equivalent to the condition imposed in [27, 3.4(ii)].

It is likely that if  $G$  is  $k$ -isotropic,  $S \subset V_\infty$ , and for every  $v$  in  $S$ ,  $G(k_v)$  is topologically simply connected, then  $M(S, G)$  is isomorphic to  $\hat{\mu}(k)$ . This has been verified for many of the classical groups and some groups of exceptional types.

A variant of Moore's theorem on the "uniqueness of the reciprocity law", announced in [25], together with the results of [28, 30], can be used to compute  $M(S, G)$ , modulo 2-torsion, for all  $k$ -anisotropic  $G$ . For some results in this direction see [32].

#### 4. Projective-Simplicity of $G(k)$ and Centrality of $C(S, G)$

It has been conjectured by Kneser, Platonov and Margulis that *if  $G$  is isotropic at each nonarchimedean place, then  $G(k)$  is projectively-simple* i.e. it does not contain any proper noncentral normal subgroup, *and if it is anisotropic at some nonarchimedean place, then (as is well known,  $G$  is of type A and) any noncentral normal subgroup of  $G(k)$  is the intersection of  $G(k)$  with a normal subgroup of  $\prod_{v \in \mathcal{S}} G(k_v)$ , where  $\mathcal{S}$  is the (finite) set of nonarchimedean places of  $k$  where  $G$  is anisotropic*. This conjecture is known to hold for all  $k$ -isotropic groups except possibly for certain outer forms of type  $E_6$  of  $k$ -rank 1 which require division algebras of degree 3 for their construction. For anisotropic groups, the results are much less complete. In 1980, inspired by [22], Margulis [16] proved the above conjecture for groups of type  $A_1$ . This implies the projective-simplicity of the spin group of any quadratic form in 3 or 4 variables which is isotropic at all nonarchimedean places of  $k$ . Projective-simplicity of the spin group of any quadratic form in at least five variables was proved already in 1956 by Kneser [9] by an ingenious method. Borovoi and Chernousov have recently proved the projective-simplicity of  $G(k)$  whenever  $G$  is of absolute rank at least two and it splits over a quadratic extension of  $k$  (this class includes all groups of type B, C,  $E_7$ ,  $E_8$ ,  $F_4$  and  $G_2$ ); and now Sury and Tomanov have independently established this for  $G$  of type  $A_3$ , which is isotropic at all nonarchimedean places – this implies the projective-simplicity of  $G(k)$  for all groups  $G$  of type D (except the triality forms). But the anisotropic groups of (inner and outer type)  $A_n$  ( $n$  arbitrary) pose a serious challenge.

It may be of interest to note here that it follows from the well known result of Margulis [15] on normal subgroups of lattices<sup>2</sup> in semi-simple groups, and

---

<sup>2</sup> A lattice in a locally compact unimodular group is a discrete subgroup of finite covolume (with respect to any Haar measure).

the strong approximation property, that any noncentral normal subgroup of  $G(k)$  is of finite index in  $G(k)$  (see, for example, [23]). Moreover, it is easy to show that if  $G$  is isotropic at all the nonarchimedean places of  $k$  and has the congruence subgroup property for some  $S$ , then  $G(k)$  is projectively-simple.

Based on the results of [2, 33] on  $\mathrm{SL}_n$  and  $\mathrm{Sp}_{2n}$ , and [18], where the centrality of the  $S$ -congruence kernel was proved for all Chevalley groups of rank  $\geq 2$ , it has been conjectured that for arbitrary (simply connected)  $G$ ,  $C(S, G)$  is central in  $\widehat{G}_S$  if  $\sum_{v \in S} k_v\text{-rank}(G) \geq 2$  and  $G$  is isotropic at all nonarchimedean  $v \in S$ . Using some of the ideas of [2, 33], Vaserstein [35] showed that this conjecture holds for all classical groups of  $k$ -rank at least two. Raghunathan has proved the above conjecture for all  $k$ -isotropic groups [31]; his proof does not require any case-by-case analysis.

For the spin group of an arbitrary (not necessarily isotropic) quadratic form in at least five variables the above conjecture on the centrality of  $C(S, G)$  was proved by Kneser [10]. Refining and using his ideas, Rapinchuk and Tomanov have recently proved the conjecture for all anisotropic groups of type  $B_r(r \geq 2)$ ,  $C_r(r \geq 2)$ ,  $D_r(r \geq 5)$ ,  $E_7$ ,  $E_8$ ,  $F_4$ ,  $G_2$ , and the groups of type  ${}^2A_r(r \geq 3)$  which split over a quadratic extension of  $k$ . The question of centrality of the  $S$ -congruence kernel for anisotropic groups of type  $A$  is a very interesting open problem—its solution may require new insight into the structure and geometry of central division algebras over global fields.

## 5. The Hasse Principle and Tamagawa Number

If  $k$  is a global function field, then the Galois cohomology  $H^1(k, G)$  is trivial (this was proved by Harder). On the other hand, if  $k$  is a number field, it has been known for quite some time that the “Hasse principle” i.e., the assertion that the natural morphism

$$H^1(k, G) \rightarrow \prod_{v \in V_\infty} H^1(k_v, G)$$

is injective, holds for all (simply connected)  $G$  of type other than  $E_8$ . The Hasse principle has now been verified for groups of type  $E_8$  by Chernousov [5].

If  $k$  is number field, let  $D_k$  be the absolute value of the discriminant of  $k/\mathbb{Q}$  and if  $k$  is a global function field, let  $q_k$  be the cardinality of its field of constants,  $g_k$  the genus of  $k$  and  $D_k = q_k^{2g_k - 2}$ .

Let  $\omega$  be an invariant exterior form on  $G$ , of maximal degree, defined over  $k$ . Then for each place  $v$ , the form  $\omega_v$ , and the normalized Haar measure on  $k_v$ , determine a Haar measure on  $G(k_v)$  which we denote by  $\omega_v$ .

Let  $P = (P_v)_{v \in V_f}$  be a fixed coherent collection of parahoric subgroups: for each  $v \in V_f$ ,  $P_v$  is a parahoric subgroup of  $G(k_v)$  such that the product  $\prod_{v \in V_\infty} G(k_v) \cdot \prod_{v \in V_f} P_v$  is an open subgroup of  $G(A)$ . (Recall that a subgroup of  $G(k_v)$  is said to be an *Iwahori subgroup* if it is the normalizer of a maximal pro- $p_v$  subgroup of  $G(k_v)$  or, equivalently, it is the stabilizer of a chamber (i.e. a simplex of maximal dimension) in the Bruhat-Tits building of  $G(k_v)$ . Any subgroup containing an Iwahori subgroup is called a *parahoric subgroup*.) It is known that,

as  $G$  is semi-simple, the product  $\prod \omega_v(P_v)$  is absolutely convergent and so there is a Haar measure  $\mu$  on  $G(A)$  which on the open subgroup  $\prod_{v \in V_\infty} G(k_v) \cdot \prod_{v \in V_f} P_v$  coincides with the measure  $D_k^{-\frac{1}{2} \dim G} \prod_{v \in V_\infty} \omega_v \cdot \prod_{v \in V_f} \omega_v|_{P_v}$ . It is obvious from the product formula (i.e.  $\prod_v |x|_v = 1$  for  $x \in k^\times$ ) that the measure  $\mu$  is independent of the  $k$ -form  $\omega$  and it is called the *Tamagawa measure*. The *Tamagawa number* of  $G/k$ , to be denoted  $\tau_k(G)$ , is the positive real number  $\mu(G(A)/G(k))$ . It was conjectured by Weil that for all (simply connected absolutely almost simple)  $G$ ,  $\tau_k(G) = 1$ . This conjecture has recently been proved by Kottwitz, over number fields, without any case-by-case considerations (see [11], and also [26, 3.3]). Using Arthur's trace formula and the Hasse principle, he has in fact shown that if  $k$  is a number field and  $\mathcal{G}$  is the unique quasi-split *inner*  $k$ -form of  $G$ ,  $\tau_k(G) = \tau_k(\mathcal{G})$ ; this result was conjectured by Langlands. Now since the Tamagawa number of any simply connected quasi-split group is 1 [12, 13], Weil's conjecture follows.

Weil's conjecture remains unproven for groups defined over global function fields. It is still unknown, for example, if over such a  $k$ , the Tamagawa number of every outer  $k$ -form of type A is 1.

## 6. Covolumes of $S$ -Arithmetic Subgroups, Bound for Class Numbers and the Finiteness Theorems

We shall now describe a formula for the covolume of  $S$ -arithmetic subgroups with respect to a natural Haar measure on  $G_S$ . We begin by describing a natural Haar measure  $\mu_v$  on  $G(k_v)$  for any place  $v$  of  $k$ . For a nonarchimedean place  $v$  of  $k$ , let  $\mu_v$  be the Tits measure on  $G(k_v)$  i.e. the Haar measure with respect to which the volume of any Iwahori subgroup of  $G(k_v)$  is 1. If  $v$  is archimedean, then  $k_v$  is either  $\mathbf{R}$  or  $\mathbf{C}$  and  $\mu_v$  is the Haar measure on  $G(k_v)$  such that, in the induced measure, any maximal compact subgroup of  $R_{k_v/\mathbf{R}}(G)(\mathbf{C})$  has volume 1. Now on  $G_S = \prod_{v \in S} G(k_v)$  we take the product measure  $\mu_S := \prod_{v \in S} \mu_v$ .

Let  $\mathcal{G}$  be the unique quasi-split *inner*  $k$ -form of  $G$ . For each nonarchimedean place  $v$ , we fix a parahoric subgroup  $\mathcal{P}_v$  of  $\mathcal{G}(k_v)$  of maximal volume such that  $\prod_{v \in V_\infty} \mathcal{G}(k_v) \cdot \prod_{v \in V_f} \mathcal{P}_v$  is an open subgroup of  $\mathcal{G}(A)$ .

As in Section 5, let  $P = (P_v)_{v \in V_f}$  be a fixed coherent collection of parahoric subgroups. Let  $S$  be a finite set of places containing  $V_\infty$  and let  $\Lambda = G(k) \cap \prod_{v \notin S} P_v$ . In its natural embedding in  $G_S$ ,  $\Lambda$  is an  $S$ -arithmetic subgroup. Let  $G_v$  denote the smooth affine  $\mathfrak{o}_v$ -group scheme associated with the parahoric subgroup  $P_v$  by the Bruhat-Tits theory ([34, 3.4]). Let  $\bar{G}_v := G_v \otimes_{\mathfrak{o}_v} \mathfrak{f}_v$  be the reduction mod  $\mathfrak{p}_v$  of  $G_v$ . Let  $\bar{T}_v$  be a maximal  $\mathfrak{f}_v$ -torus of  $\bar{G}_v$  containing a maximal  $\mathfrak{f}_v$ -split torus and  $\bar{M}_v$  be the maximal reductive  $\mathfrak{f}_v$ -subgroup containing  $\bar{T}_v$ . Note that  $\bar{M}_v$  depends on the choice of  $P_v$ . Let  $\bar{\mathcal{G}}_v, \bar{\mathcal{T}}_v$  and  $\bar{\mathcal{M}}_v$  be similarly defined  $\mathfrak{f}_v$ -groups associated with  $\mathcal{G}$  and the parahoric subgroup  $\mathcal{P}_v$ .

If  $\mathcal{G}/k$  is not a triality form of type  ${}^6D_4$ , let  $\ell$  be the smallest extension of  $k$  over which  $\mathcal{G}$  splits. If  $\mathcal{G}/k$  is of type  ${}^6D_4$ , let  $\ell$  be a fixed extension of  $k$  of degree 3 contained in the Galois extension of degree 6 over which  $\mathcal{G}$  splits. Let  $D_\ell$  be the absolute value of the discriminant of  $\ell/\mathbf{Q}$  if  $k$  is a number field and

$D_\ell = q_\ell^{2g_\ell - 2}$  if  $k$  is a global function field, where  $q_\ell$  is the cardinality of the finite field of constants in  $\ell$  and  $g_\ell$  is the genus of  $\ell$ .

*The integer  $s(\mathcal{G})$ :* If  $\mathcal{G}$  splits over  $k$ , let  $s(\mathcal{G}) = 0$ . If  $\mathcal{G}$  is a  $k$ -form of type  ${}^2A_r$ , with  $r$  even, let  $s(\mathcal{G}) = \frac{1}{2}r(r+3)$ ; if  $\mathcal{G}$  is a  $k$ -form of type  ${}^2D_r$  ( $r$  arbitrary) or  ${}^2E_6$ , let  $s(\mathcal{G}) = \frac{1}{2}(r-1)(r+2), 2r-1$  or  $26$  respectively. If  $\mathcal{G}$  is a triality form of type  ${}^3D_4$  or  ${}^6D_4$ , then let  $s(\mathcal{G}) = 7$ .

The following theorem provides a “computable” formula for the volume of  $S$ -arithmetic quotients of  $G_S$ . It is proved in [26].

**Theorem 4.** *Let  $m_1 \leq \dots \leq m_r$  be the exponents of the Weyl group of the absolute root system of  $G$ . Then*

$$\mu_S(G_S/A) = D_k^{\frac{1}{2} \dim G} (D_\ell / D_k^{[\ell:k]})^{\frac{1}{2}s(\mathcal{G})} \left( \prod_{v \in V_\infty} \left| \prod_{i=1}^r \frac{m_i!}{(2\pi)^{m_i+1}} \right|_v \right) \tau_k(G)\mathcal{E}(P) ;$$

where

$$\mathcal{E}(P) = \prod_{v \in S_f} \frac{q_v^{(r_v + \dim \overline{M}_v)/2}}{\# \overline{T}_v(\mathfrak{f}_v)} \cdot \prod_{v \notin S} \frac{q_v^{(\dim \overline{M}_v + \dim \overline{M}_v)/2}}{\# \overline{M}_v(\mathfrak{f}_v)} ,$$

$S_f = S \cap V_f$ , and for  $v \in V_f$ ,  $r_v (= \dim \overline{T}_v)$  is the rank of  $G$  over the maximal unramified extension of  $k_v$ .

The results involved in the proof of this theorem provide the following lower bound for the class number of simply connected anisotropic groups (see [26, Theorem 4.3]).

**Theorem 5.** *Assume that  $G$  is anisotropic over  $k$  and moreover  $G_\infty := \prod_{v \in V_\infty} G(k_v)$  is compact. Then the class number  $\#(G_\infty \prod_{v \in V_f} P_v \backslash G(A)/G(k))$  of  $G/k$  with respect to  $P$  is at least*

$$D_k^{\frac{1}{2} \dim G} (D_\ell / D_k^{[\ell:k]})^{\frac{1}{2}s(\mathcal{G})} \left( \prod_{v \in V_\infty} \left| \prod_{i=1}^r \frac{m_i!}{(2\pi)^{m_i+1}} \right|_v \right) \tau_k(G)\zeta(P) ;$$

where

$$\zeta(P) = \prod_{v \in V_f} \frac{q_v^{(\dim \overline{M}_v + \dim \overline{M}_v)/2}}{\# \overline{M}_v(\mathfrak{f}_v)} .$$

In [3, §7] this theorem is used to prove the following finiteness theorem.

**Theorem 6.** *Given a positive integer  $c$ , let  $\mathcal{C}_c$  be the set of pairs  $(k, G)$  consisting of a number field  $k$  and a connected, simply connected absolutely almost simple group  $G$  such that  $G$  is anisotropic over  $k$ ,  $G_\infty := \prod_{v \in V_\infty} G(k_v)$  is compact, and the class number  $\#(G_\infty \prod_{v \in V_f} P_v \backslash G(A)/G(k))$  of  $G/k$  with respect to some coherent collection of parahoric subgroups  $(P_v)_{v \in V_f}$  is less than  $c$ . Then (up to natural equivalence)  $\mathcal{C}_c$  is finite.*

The formula for the volume of  $S$ -arithmetic quotients given above and certain number theoretic estimates have been used in [3] to prove that in characteristic zero, there are only finitely many distinct  $S$ -arithmetic subgroups  $\Gamma$  of covolume  $\leq c$ , where  $c$  is a given positive number. Also, there are only finitely many  $S$ -arithmetic  $\Gamma$  with  $0 \neq |\chi(\Gamma)| < c$ , where  $\chi(\Gamma)$  is the Euler-Poincaré characteristic of  $\Gamma$  in the sense of C. T. C. Wall. For precise results, see [3].

## 7. Super-Rigidity and Arithmeticity of Lattices

According to a celebrated theorem of Margulis (announced at the ICM held in 1974), *irreducible* lattices<sup>3</sup> in real semi-simple groups of  $\mathbf{R}$ -rank  $> 1$  are super-rigid<sup>4</sup>. It follows rather easily from this that such lattices are arithmetic. On the other hand, it has been known for almost thirty years that the groups  $\mathrm{SO}(n, 1)$  contain non-arithmetic lattices for  $n \leq 5$ . In 1986, Gromov and Piatetski-Shapiro, employing a nice geometric construction, showed that for each  $n$ ,  $\mathrm{SO}(n, 1)$  contains plenty of non-arithmetic lattices([8]). Mostow has constructed non-arithmetic lattices in  $\mathrm{SU}(2, 1)$  and  $\mathrm{SU}(3, 1)$ ; however, whether lattices in  $\mathrm{SU}(n, 1)$  are all arithmetic if  $n$  is sufficiently large is still an open question.

Corlette [6] has now established the super-rigidity of real representations of *cocompact* discrete subgroups in the remaining semi-simple groups of  $\mathbf{R}$ -rank 1, namely the groups  $\mathrm{Sp}(n, 1)$  and the  $\mathbf{R}$ -rank 1 form of type  $F_4$ , using his basic theorem on the existence of a harmonic map in any given homotopy class of maps from a compact riemannian manifold into a locally symmetric space; and just a few weeks ago I have learnt that Gromov and Schoen have proved that any representation of such a discrete cocompact subgroup over a  $p$ -adic field is bounded by developing an analogue of the theory of harmonic maps for maps from a riemannian manifold into a Bruhat-Tits building. Now, as in the case of groups of  $\mathbf{R}$ -rank  $> 1$ , arithmeticity of cocompact discrete subgroups of the groups  $\mathrm{Sp}(n, 1)$  and the  $\mathbf{R}$ -rank 1 form of type  $F_4$  follows.

*Acknowledgements.* I thank the University of Michigan for inviting me to deliver the Fall 1989 Keeler lectures, and Mr. M.S. Keeler for instituting this lectureship. This address is in part based on my Keeler lectures. I also thank the Japan Association for Mathematical Sciences for financial support to attend the ICM.

## References

1. Bak, A., Rehmann, U.: The congruence subgroup and metaplectic problems for  $\mathrm{SL}_{n \geq 2}$  of division algebras. *J. Alg.* **78** (1982) 475–547
2. Bass, H., Milnor, J., Serre, J.-P.: Solution of the congruence subgroup problem for  $\mathrm{SL}_n$  and  $\mathrm{Sp}_{2n}$ . *Publ. Math. IHES* **33** (1967) 59–137

<sup>3</sup> A lattice in a semi-simple linear analytic group is said to be *irreducible* if no subgroup of it of finite index is a direct product of two infinite normal subgroups.

<sup>4</sup> For a detailed proof see [17]. Super-rigidity and arithmeticity in arbitrary characteristic has been proved in [36].

3. Borel, A., Prasad, G.: Finiteness theorems for discrete subgroups of bounded covolume in semi-simple groups. *Publ. Math. IHES* **69** (1989) 119–171; Addendum: *ibid* **71** (1990)
4. Chase, S., Waterhouse, W.: Moore's theorem on uniqueness of reciprocity laws. *Invent. math.* **16** (1972) 267–270
5. Chernousov, V.: On the Hasse principle for groups of type  $E_8$ . *Soviet Math. Dokl.* **39** (1989)
6. Corlette, K.: Archimedean super-rigidity and hyperbolic geometry (preprint)
7. Deodhar, V.: On central extensions of rational points of algebraic groups. *Amer. J. Math.* **100** (1978) 303–386
8. Gromov, M., Piatetski-Shapiro, I.: Non-arithmetic groups in Lobachevsky spaces. *Publ. Math. IHES* **66** (1988) 93–103
9. Kneser, M.: Orthogonale Gruppen über algebraischen Zahlkörpern. *J. Reine Angew. Math.* **196** (1956) 213–220
10. Kneser, M.: Normalteiler ganzzahliger Spingruppen. *J. Reine Angew. Math.* **311/312** (1979) 191–214
11. Kottwitz, R.: Tamagawa numbers. *Ann. Math.* **127** (1988) 629–646
12. Langlands, R.: The volume of the fundamental domain. *Proc. AMS Symp. Pure Math.* **9** (1966) 143–148
13. Lai, K.: Tamagawa number of reductive algebraic groups. *Compos. Math.* **41** (1980) 153–188
14. Margulis, G.: Cobounded subgroups of algebraic groups over local fields. *Funct. Anal. Appl.* **11** (1977) 119–128
15. Margulis, G.: Finiteness of quotient groups of discrete subgroups. *Funct. Anal. Appl.* **13** (1979) 178–187
16. Margulis, G.: On the multiplicative group of a quaternion algebra over a global field. *Soviet Math. Dokl.* **21** (1980) 780–784
17. Margulis, G.: Arithmeticity of the irreducible lattices in the semi-simple groups of rank greater than 1. *Invent. math.* **76** (1984) 93–120
18. Matsumoto, H.: Sur les sous-groupes arithmétiques des groupes semi-simples déployés. *Ann. Sci. École Norm. Sup., 4<sup>e</sup> sér.* **2** (1969) 1–62
19. Moore, C.: Group extensions of  $p$ -adic and adelic linear groups. *Publ. Math. IHES* **35** (1968) 157–222
20. Moore, C.: Group extensions and cohomology for locally compact groups, III and IV. *Transactions AMS* **221** (1976) 1–58
21. Platonov, V.: The problem of strong approximation and the Kneser-Tits conjecture. *Math. USSR Izv.* **3** (1969) 1139–1147; Addendum: *ibid* **4** (1970) 784–786
22. Platonov, V., Rapinchuk, A.: On the group of rational points of three-dimensional groups. *Soviet Math. Dokl.* **20** (1979) 693–697
23. Prasad, G.: Strong approximation. *Ann. Math.* **105** (1977) 553–572
24. Prasad, G.: Elementary proof of a theorem of Bruhat-Tits and Rousseau. *Bull. Soc. Math. France* **110** (1982) 197–202
25. Prasad, G.: A variant of a theorem of Calvin Moore. *C. R. Acad. Sci. (Paris) Sér. I* **302** (1982) 405–408
26. Prasad, G.: Volumes of  $S$ -arithmetic quotients of semi-simple groups. *Publ. Math. IHES* **69** (1989) 91–117
27. Prasad, G., Raghunathan, M.S.: On the congruence subgroup problem: Determination of the “metaplectic kernel”. *Invent. math.* **71** (1983) 21–42
28. Prasad, G., Raghunathan, M.S.: Topological central extensions of semi-simple groups over local fields. *Ann. Math.* **119** (1984) 143–268
29. Prasad, G., Raghunathan, M.S.: On the Kneser-Tits problem. *Comment. Math. Helv.* **60** (1985) 107–121

30. Prasad, G., Raghunathan, M.S.: Topological central extensions of  $\mathrm{SL}_1(D)$ . *Invent. math.* **92** (1988) 645–689
31. Raghunathan, M.S.: On the congruence subgroup problem. *Publ. Math. IHES* **46** (1976) 107–161; second part: *Invent. math.* **85** (1986) 73–117
32. Rapinchuk, A.: Multiplicative arithmetic of division algebras over number fields and the metaplectic problem. *Math. USSR Izv.* **31** (1988) 349–379
33. Serre, J.-P.: Le problème des groupes de congruence pour  $\mathrm{SL}_2$ . *Ann. Math.* **92** (1970) 489–527
34. Tits, J.: Reductive groups over local fields. *Proc. AMS Symp. Pure Math.* **33** (1979), part 1, 29–69
35. Vaserstein, L.: The structure of classical arithmetic groups of rank greater than one. *Math. USSR Sb.* **20** (1973) 465–492
36. Venkataramana, T.: On super-rigidity and arithmeticity of lattices in semi-simple groups over local fields of arbitrary characteristic. *Invent. math.* **92** (1988) 255–306
37. Wigner, D.: Algebraic cohomology of topological groups. *Transactions AMS* **178** (1973) 83–93

# Poles of Standard $L$ Functions

*Stephen Rallis*

Department of Mathematics, The Ohio State University, Columbus, OH 43210, USA

## Introduction

The Rankin-Selberg method in the theory of  $L$  functions gives explicit integral representations of certain  $L$  functions of automorphic representations of reductive algebraic groups. This allows one to determine the analytic continuation and functional equations of such  $L$  functions. Moreover the poles of these  $L$  functions can be determined explicitly. The major applications of this method include (i) the determination of how much of the complementary series contributes to cuspidal representations of  $\mathrm{GL}_2$  and (ii) the strong multiplicity one Theorem and the classification of automorphic representations of  $\mathrm{GL}_n$ .

In recent years Rankin-Selberg integral representations have been found for several new classes of  $L$  functions [PS-R-II, PS-R-S]. We are concerned here with the standard  $L$  functions of the classical groups. The method of doubling a classical group (which is essentially a compactification of the group) is used in this case. The doubling method has the extra advantage that the special values of such  $L$  functions at integral points can be related to certain  $\theta$  integrals arising from specific dual reductive pairs.

This leads to our second theme. Namely as a generalization of the classical Siegel-Weil identity, there exist, for dual reductive pairs, identities between certain regularized  $\theta$  integrals and certain special values of Siegel type Eisenstein series. The point of such identities is to give criteria about the existence of a pole of the  $L$  functions mentioned above at a specific value in terms of the nonvanishing of a certain  $\theta$  lift between groups of a given dual reductive pair.

The work in this lecture represents joint work (as specified below) with Ilya Piatetski Shapiro and Steve Kudla. For a survey of earlier related works, and a more extensive bibliography of the field, see [G-S].

## §0. Notation

(1) Let  $G$  be a reductive group over  $k$ , a number field. Let  $G_v$  the associated group at the place  $v$  of  $k$ . Let  $G(\mathbb{A})$  be the corresponding adelic group.

Let  $\mathbb{K} = \prod_v K_v$  be a good maximal compact group where  $K_v$  is a special maximal compact subgroup of  $G_v$  (for all  $v$ ).

Let  $\Pi = \bigotimes_v \Pi_v$  be an irreducible automorphic cuspidal representation of  $G(\mathbb{A})$ . Define  $S_\Pi$  = the set of places of  $k$  consisting of the Archimedean places together with all finite places at which  $\Pi_v$  is not spherical. For  $v \notin S_\Pi$  let  $D_v(\Pi_v)$  be the conjugacy class in the  $L$  group of  $G_v$  associated to  $\Pi_v$  via the Satake isomorphism. Then for a representation  $\varrho$  of the  $L$  group  ${}^L G$  let

$$L_v(s, \Pi_v, \varrho) = \det(1 - \varrho(D_v)q_v^{-s})^{-1}$$

be the usual Langlands Euler factor with  $s \in \mathbb{C}$  and  $q_v$  = cardinality of residue field of  $k_v$ . Moreover let

$$L_S(s, \Pi, \varrho) = \prod_{v \notin S_\Pi} L_v(s, \Pi_v, \varrho)$$

be the corresponding global restricted  $L$  function.

We let  $\zeta_v(s)$  be the usual zeta function associated to a local field  $k_v$ .

(2) We let  $\mathrm{Sp}_n \times \mathrm{O}(Q)$  be a dual reductive pair, where  $\mathrm{Sp}_n$  is the symplectic group of  $2n \times 2n$  matrices and  $\mathrm{O}(Q)$  is the orthogonal group of the quadratic form  $Q$  with  $\dim Q$  even. The oscillator representation of this pair is denoted by  $\omega_{Q,\psi}$  ( $\psi$  an additive character) and is realized on the space  $S(M_{mn}(k))$  ( $m = \dim Q$ ) (See [W]). The construction works locally and globally (where the local field  $k$  is replaced by  $\mathbb{A}_k = \mathbb{A}$  = adeles of the number field  $k$ ). Following the work of [W] it is possible to define the space of  $\theta$  kernels

$$\theta_\varphi(x, y) = \sum_{\xi \in M_{mn}(k)} \omega_{Q,\psi}(x, y)(\varphi)(\xi)$$

where  $\varphi \in S(M_{mn}(\mathbb{A}))$  and  $(x, y) \in \mathrm{Sp}_n \times \mathrm{O}(Q)(\mathbb{A})$ . The kernel  $\theta_\varphi$  is left invariant by  $\mathrm{Sp}_n \times \mathrm{O}(Q)(k)$  and slowly increasing on  $\mathrm{Sp}_n \times \mathrm{O}(Q)(k) \backslash \mathrm{Sp}_n \times \mathrm{O}(Q)(\mathbb{A})$ . In particular this allows one to define for any  $f$ , which is rapidly decreasing on  $\mathrm{O}(Q)(k) \backslash \mathrm{O}(Q)(\mathbb{A})$ , the  $\theta$  integral  $\theta_\varphi(f)$  as

$$\int_{\mathrm{O}(Q)(k) \backslash \mathrm{O}(Q)(\mathbb{A})} \theta_\varphi(x, y) f(y) dy .$$

A similar  $\theta_\varphi(F)$  can be defined for  $F$  rapidly decreasing on  $\mathrm{Sp}_n(k) \backslash \mathrm{Sp}_n(\mathbb{A})$ .

(3) Given a reductive group  $G$  we let  $\mathcal{A}(G(\mathbb{A}))$  be the space of slowly increasing automorphic functions on  $G(k) \backslash G(\mathbb{A})$  as given in [B-J]. Let  $L^2_{\mathrm{cusp}}(G(\mathbb{A}))$  and  $L^2_{\mathrm{res}}(G(\mathbb{A}))$  be the space of smooth cusp forms and smooth residual forms in  $L^2(G(k) \backslash G(\mathbb{A}))$ . A function  $f$  belonging to one of these spaces is smooth if  $f$  is fixed by a compact open subgroup  $K'$  of  $G(\mathbb{A}_{\mathrm{fin}})$  and  $f$  is a  $C^\infty$  vector relative to  $G_\infty$  component of  $G(\mathbb{A})$ .

## §1. The Doubling Method

We describe a setup generalizing the construction of [PS-R-II]. We let  $K$  be a 2 dimensional semisimple commutative algebra over a number field  $k$ . Then  $K$  is either a quadratic extension of  $k$  or  $K = k \oplus k$ . We let  $V$  be a finite dimensional  $k$  vector space provided with a nondegenerate  $\varepsilon$  symmetric bilinear form  $\langle , \rangle (\varepsilon = \pm 1)$ . We consider the extension of scalars functor  $\sim$  applied to  $V, \langle , \rangle$ . Then  $\langle , \rangle \sim = \langle , \rangle \otimes_k K$  is a  $K$  valued bilinear form on the space  $V \otimes_k K = W$ . If  $\text{tr}_{K/k}$  is the canonical trace form on  $K$  then there exists an element  $\xi \in K^\times$  so that the form  $\langle\langle , \rangle\rangle = \text{tr}_{K/k}(\xi \langle , \rangle \sim)$  is a totally split  $\varepsilon$  symmetric bilinear form on  $W$  over  $k$  (if  $K$  is a quadratic field  $k(\sqrt{d})$  then we let  $\xi = \sqrt{d}$  and if  $K = k \oplus k$  we let  $\xi = (1, -1)$ ).

We let  $G = U(W, \langle\langle , \rangle\rangle)$  be the group of  $k$  linear isometries of  $W, \langle\langle , \rangle\rangle$ . If  $2\ell = \dim_k(W)$ , let  $W = W_{2\ell}$ . We let  $H = U_K(W, \langle , \rangle \sim)$  be the group of  $K$  linear maps of  $W$  which preserve  $\langle , \rangle \sim$ . Then  $H \subseteq G$ .

Moreover if  $K = k \oplus k$  then  $H \cong H_1 \times H_1$  where  $H_1$  is  $U(V, \langle , \rangle)$ , the  $k$  isometry group of  $V, \langle , \rangle$ .

We let  $X_r =$  variety of  $r$ -dimensional  $\langle\langle , \rangle\rangle$  isotropic subspaces of  $W$  (as a  $k$  space). Then there is a parabolic subgroup  $P_r$  of  $G$  so that  $X_r \cong P_r \backslash G$ . The set of double cosets  $P_r \backslash G / H$  in general is not finite. However it is finite precisely when  $\varepsilon = -1$  (for all  $r$ ) and when  $\varepsilon = 1$  with  $r = \frac{1}{2} \dim_k(W) = \dim_k(V)$ .

We are interested in determining the orbits of  $H$  in  $P_r \backslash G$  which are negligible. Namely an  $H$  orbit  $\mathcal{O}$  in  $X_r$  is negligible if the stabilizer in  $H$  of a subspace  $T \in \mathcal{O}$  contains as a normal subgroup the unipotent radical of a proper parabolic subgroup of  $H$ . Then in the case  $P_r \backslash G / H$  is finite the number of nonnegligible orbits is exactly one in the cases  $r = 2, 3, \dim V - 1, \dim V (\varepsilon = -1)$  and  $r = \dim V (\varepsilon = 1)$  ([Rab]).

Since the form  $\langle\langle , \rangle\rangle$  is totally split the group  $P_r = M_r U_r$  where  $U_r$  is the unipotent radical of  $P_r$  and  $M_r \cong \text{GL}_r(k) \times U(W_{2\ell-2r}, \langle\langle , \rangle\rangle_{2\ell-2r})$  with  $\langle\langle , \rangle\rangle_{2\ell-2r}$  a totally split  $\varepsilon$  symmetric bilinear form on  $W_{2\ell-2r}$  with  $W_{2\ell-2r} \subseteq W_{2\ell}$ . Then we consider the one dimensional character  $\delta_{P_r} : P_r \rightarrow k^\times$  defined by  $\delta_{P_r}(g) = \det(\text{Ad}_{U_r}(g))$  with  $g = mu$  relative to the decomposition above.

For each place  $v$  of  $k$  we define the quasi-character  $\chi_{s,v}$  of  $(M_r)_v$  by

$$\chi_{s,v}(m) = |\delta_{P_r}(m)|^{\kappa_r s}$$

with  $\kappa_r^{-1} = 2\ell - r + 1$  ( $2\ell - r - 1$  resp.) if  $\varepsilon = -1$  ( $\varepsilon = 1$  resp.)

Then define the induced representation

$$I_v(s) = \text{Ind}_{(P_r)_v}^{G_v}(\chi_{s,v})$$

where the induction is normalized. The global induced representation is given by

$$I(s) = \text{Ind}_{P(\mathbb{A})}^{G(\mathbb{A})}(\bigotimes_v \chi_{s,v}) \simeq \bigotimes_v I_v(s).$$

Starting with a section  $\phi(\cdot, s) \in I(s)$  we form the Eisenstein series

$$E_\ell^r(g, s, \phi) = \sum_{\gamma \in P_r(k) \backslash G(k)} \phi(\gamma g, s).$$

For a fixed good maximal compact subgroup  $\mathbb{K} = \prod_v K_v$  of  $G(\mathbb{A})$  we say  $\phi(\cdot, s) \in I(s)$  is standard if the restriction of  $\phi(\cdot, s)$  to  $\mathbb{K}$  is independent of  $s$  and  $\phi(\cdot, s)$  is  $K_\infty$  finite for each  $K_\infty$ ,  $\infty$  an Archimedean prime.

We let  $\Pi = \bigotimes_v \Pi_v$  be an irreducible automorphic cuspidal representation of the group  $H(\mathbb{A})$ . We let  $f_\Pi \in \Pi$ .

Then we form the Rankin-Selberg integral

$$Z_r(f_\Pi, \phi(\cdot, s)) = \int_{H(k) \backslash H(\mathbb{A})} f_\Pi(h) E_\ell^r(h, s, \phi) dh$$

relative to some choice of suitable Tamagawa measure  $dh$  on  $H(k) \backslash H(\mathbb{A})$  and  $\phi$  standard.

**Lemma 1.1.** *In the case  $P_r \backslash G/H$  admits one nonnegligible orbit (see conditions above)*

$$Z_r(f_\Pi, \phi(\cdot, s)) = \int_{H_\gamma^0(\mathbb{A}) \backslash H(\mathbb{A})} \phi(\gamma g, s) \left( \int_{H_\gamma(k) \backslash H_\gamma^0(\mathbb{A})} f_\Pi(hg) dh \right) dg$$

where  $H_\gamma(k)$  is the stabilizer of the point  $P_r \gamma$  in the unique nonnegligible orbit of  $H$  in  $P_r \backslash G$  and  $H_\gamma^0(\mathbb{A})$  is the group of norm one elements of  $H_\gamma(\mathbb{A})$ .

*Proof.* The proof is the usual method of unwinding the series expansion of  $E^r$  relative to the  $H(k)$  orbits in  $P_r \backslash G$ . We note that a negligible orbit contributes zero since (i) the stabilizer contains a normal subgroup which is the unipotent radical of a proper parabolic of  $H$  and (ii) the form  $f_\Pi$  is cuspidal. Thus we are left with one term as above (with the hypotheses of the Lemma).  $\square$

*Remark.* 1.1. We consider now the case where  $\dim(V) = r$ . In fact in these examples

$$H_\gamma(k) \cong \begin{cases} \{(g, g) | g \in H_1\} \cong H_1 & \text{if } K = k \oplus k \\ \text{the } k \text{ rational points of } H \simeq H_1 & \text{if } K \text{ is quadratic over } k. \end{cases}$$

The inner integral in Lemma 1.1 can be calculated. We discuss the two cases separately.

(i) If  $K = k \oplus k$  and  $\Pi = \Pi_1 \otimes \Pi_2, \Pi_i$  automorphic cuspidal irreducible representation of  $H_1(\mathbb{A})$ , the inner integral is

$$\int_{H_1(k) \backslash H_1(\mathbb{A})} f_{\Pi_1}(hg_1) f_{\Pi_2}(h) dh.$$

Such a term is zero unless  $\Pi_1 \cong \Pi_2^\sim =$  the contragredient of  $\Pi_2$ . Moreover, the integral (with data defining  $f_{\Pi_1}$  and  $f_{\Pi_2}$  suitably normalized) then becomes a matrix coefficient of  $\langle f_{\Pi_1} * g_1 | f_{\Pi_2^\sim} \rangle$  of  $\Pi_1$ .

(ii) If  $K$  is a field, the inner integral is

$$\int_{H_1(k) \backslash H_1(\mathbb{A})} f_H(hg_1) dh.$$

This represents the period of  $f_H$  over the subgroup  $H_1 \subset H$ . A fundamental question is to determine for which  $\Pi$  such an integral is nonvanishing. We note that a similar problem exists in the doubling for the  $GL_n$  case ([PS-R-II]). In fact in that case the set of automorphic representation  $\Pi$  which have nonvanishing period should be the set of  $\Pi$  which arise by quadratic base change from an appropriate unitary group (care must be taken in that case for the correct condition on central characters).

## §2. Factorizability Properties

The point of expressing  $Z_r$  as an integral in Lemma 1.1 above is to determine whether the inner integral can be factorized as an infinite product (Euler Product) of local factors.

Indeed in the case  $K = k \oplus k$  the general matrix coefficient of  $\Pi_1$  factors as

$$\langle f_{\Pi_1} * g | f_{\Pi_1^\sim} \rangle = \prod_v \langle \xi_v * g_v | \xi_v^\sim \rangle_{(\Pi_1)_v}$$

provided there exists an embedding  $\Pi \xrightarrow{i_\Pi} L^2_{\text{cusp}}(H_1(k) \backslash H_1(\mathbb{A}))$  (which is  $H_1(\mathbb{A})$  intertwining) in such a way that  $i_{\Pi_1}(\otimes \xi_v) = f_{\Pi_1}, i_{\Pi_1^\sim}(\otimes \xi_v^\sim) = f_{\Pi_1^\sim}$ . The point here is that there is for each prime  $v \in k$  exactly one  $H_1 \cong H_v(k_v)$  invariant form on the space  $(\Pi_1)_v \otimes (\Pi_1^\sim)_v$  for each irreducible admissible representation  $\Pi_1$  (Schur's Lemma).

This example illustrates the very general principle of uniqueness that explains when a global Rankin integral (such as the general  $Z_r$  above) has an Euler product (independent of making the usual unwindings of the integral). Indeed in the example of  $Z_r$  we look at the following local problem. The global integral defining  $Z_r$  determines for each place  $v$  a bilinear functional on  $(\Pi_1)_v \otimes I_v(s)$  which is  $H_v$  invariant. Thus we consider the space  $\text{Hom}_{H_v}((\Pi_1)_v \otimes I_v(s), 1)$  for all  $s$ .

If the dimension of such a space is at most 1 (for all  $s$ ) and if this holds for all primes  $v$ , then  $Z_r()$  can be factorized as an infinite product of local factors. This is just the simple principle that local multiplicity one (for all primes  $v$ ) implies global multiplicity one. The main problem here is that in order to express  $Z_r$  as such a factorizable product we need to know for some a priori reason that in fact the global Rankin integral is nonvanishing. However in the case  $Z_r$  given in Lemma 1.1 the nonvanishing property is determined by the nonvanishing of a certain period (see (ii) in Remark 1.1). What is known is that the principle of

local uniqueness implies also the existence of the local functional equation and thus the determination of a local  $\varepsilon$  factor of a representation (see [PS-R-I] for the consideration of these matters in the doubling case).

*Remark.* 2.1. A typical Rankin-Selberg integral leads to the following type of uniqueness question of models. Namely let  $M \subseteq N$  be reductive groups. Let  $\Pi_1$  be an irreducible admissible representation of  $M(k_v)$  and  $\text{Ind}_{P_{N_v}}^{N_v}(\sigma_v)$  an induced representation from a standard parabolic  $P_{N_v}$  of  $N_v$  with  $\sigma_v$ , a character on  $P_{N_v}/N_v$ . Then we consider the space  $Y_v = \text{Hom}_{M_v}(\Pi_1 \otimes \text{Ind}_{P_{N_v}}^{N_v}(\sigma_v), 1)$ . In the particular case  $\Pi_1 = \text{Ind}_{B_v}^{M_v}(\chi_v)$  (where  $B_v$  is a Borel subgroup of  $M_v$  and  $\chi_v$  any quasicharacter on  $B_v$ ), then the determination of  $Y_v$  reduces by the usual Bruhat theory to finding certain quasi-invariant distributions on  $N_v$  relative to the left action of  $B_v$  and the right action of  $P_{N_v}$ . The importance of using such a  $\Pi_1$  is that “generically”  $\Pi_1$  is an unramified spherical representation and thus will be a local constituent of a global representation of  $M(\mathbb{A})$ . In any case the relevant local analysis is concerned with the set of double cosets  $B_v \backslash N_v / P_{N_v}$ . If the set  $B_v \backslash N_v / P_{N_v}$  is finite then we might expect that  $Y_v$  above is at most one dimensional (for “generic” values of  $\sigma_v$  and  $\chi_v$ ). We note the following analogy from the category of algebraic group representations that makes such a statement at least plausible. From [Kim] if the set  $B(\mathbb{C}) \backslash N(\mathbb{C}) / P_N(\mathbb{C})$  is finite then every finite dimensional irreducible representation of  $N(\mathbb{C})$ , whose highest weight vector is fixed by  $P_N(\mathbb{C})$  projectively, is multiplicity free when restricted to  $M(\mathbb{C})$  (with the assumption that  $N(\mathbb{C})$  and  $M(\mathbb{C})$  are connected and semisimple). This is a multiplicity one Branching Rule. It is not clear at this point whether such a multiplicity one Branching Rule determines when  $Y_v$  is one dimensional but such a possibility bears future investigation.

*Remark.* 2.2. The Rankin-Selberg integral given in [PS-R-S] is a special case of the data in Remark 2.1. Namely  $M = G_2$  and  $N = \text{SO}(7)$  and the embedding of  $G_2$  into  $\text{SO}(7)$  is given by the standard action of  $G_2$  on the space of trace zero elements of the 8 dimensional space of octonions. Here  $P_N$  is the parabolic subgroup of  $\text{SO}(7)$  which stabilizes a 2 dimensional isotropic flag. The  $L$  function represented in this example is the tensor product of  $G_2 \times \text{GL}_2$ .

In the case  $K = k \oplus k$  the zeta integral  $Z_\ell$  in Lemma 1.1 (with  $\ell = \dim V$ ) equals the product (with  $\phi(s) = \bigotimes_v \phi_v(s)$ )

$$\prod_v Z_{\ell,v}(\xi_v, \xi_v^\sim, s)$$

where  $Z_{\ell,v}(\xi_v, \xi_v^\sim, s)$  is the local zeta integral given by

$$\int_{H_1(k_v)} \langle \xi_v * g_v | \xi_v^\sim \rangle \phi_v(\gamma_v(g_v, 1), s) dg_v .$$

Then from [PS-R-I, PS-R-II] we have

**Proposition 2.1.** Let  $K = k \oplus k$  and fix an arbitrary point  $s_0 \in \mathbb{C}$ . It is possible to choose data  $\phi = \otimes \phi_v, f_{\Pi_1}$ , and  $f_{\Pi_1^\sim}$  so that

$$Z_\ell(f_{\Pi_1} \otimes f_{\Pi_1^\sim}, \phi, s) = \frac{L_S(\frac{1}{2} + s, \Pi, r) Z_\infty(s)}{b_{\ell, S}(s)}$$

where

(i)  $L_S(s, \Pi, r)$  is the restricted  $L$  function associated to the standard representation  $r$  of the  $L$  group of  $U(V, \langle , \rangle)_0$  ( $=$  the connected component of  $U(V, \langle , \rangle)$ ).

$$(ii) b_{\ell, S}(s) = \prod_{v \notin S} \left\{ \prod_{k=1}^{k=[\frac{\ell}{2}]} \zeta_v(2s + \ell + 1 - 2k) \right\} \begin{cases} \zeta_v(s + \frac{\ell+1}{2}) & \text{if } \varepsilon = -1 \\ 1 & \text{if } \varepsilon = 1, \end{cases}$$

(iii)  $Z_\infty(s)$  is a meromorphic function in  $s$  which is nonvanishing at the arbitrarily chosen point  $s_0$ . Note that the choice of data above depends on  $s_0$ .

*Remark.* 2.3. The functional equation of  $L_S(s, \Pi, r)$  is given in [PS-R-I]. In fact it is possible to define local  $L_v$  factors at the bad primes of  $\Pi$  (those  $v$  which are finite and  $\Pi_v$  is not spherical). Using then an extended definition of  $L$  we give the functional equation of the new  $L$  in [PS-R-III]. There we use the notion of local  $\varepsilon$  factors mentioned above.

*Remark.* 2.4. We note that if  $L_S(s, \Pi, r)$  admits a pole at  $s = s_0$  then  $b_{\ell, S}(s) E'_\ell(g, s, \phi)$  admits a nonzero pole at  $s = 2s_0 - 1$  for some section  $\phi$ . In fact by (iii) of Proposition 2.1 and for suitable choice of data  $b_{\ell, S}(s) Z_\ell(f_{\Pi_1} \otimes f_{\Pi_1^\sim}, \phi, s)$  has a pole at  $s = 2s_0 - 1$  and this pole comes from the pole of the normalized Eisenstein series.

### §3. Poles of Eisenstein Series

To determine information about the possible poles of  $L$  functions associated to the doubling method we require knowledge about the poles of  $E_n^n(\cdots)$ . As we shall see the residues are explicitly describable and generically square integrable.

For the remaining 2 sections we assume that the number field  $k$  is totally real.

For the discussion below we adopt a slightly more general convention. Namely we let  $G_0$  be the isometry group of the form

$$\begin{pmatrix} 0 & I_n \\ \varepsilon I_n & 0 \end{pmatrix}.$$

Thus

$$G_0 = \begin{cases} \mathrm{Sp}_n & \text{if } \varepsilon = -1 \\ \mathrm{O}(n, n) & \text{if } \varepsilon = +1. \end{cases}$$

We consider the parabolic  $P_n \cong \mathrm{GL}_n \ltimes U_n$  of  $G_0$  with  $U_n$ , the unipotent radical of  $P_n$ . We form the family of induced representations  $\mathbb{I}_n(s) = \mathrm{Ind}_{P_n(\mathbb{A})}^{G_0(\mathbb{A})} (|\det g|_{\mathbb{A}^x}^s \otimes 1_{U_n})$  (normalized induction) where  $||_{\mathbb{A}^x}$  = the usual adelic norm on  $\mathbb{A}^x$ . We form the family of Eisenstein series  $E_n^n(g, s, \phi)$  where  $\phi$  is a standard section.

**Theorem 3.1** [K-R-IV]. Let  $S$  be the set of primes  $v$  of  $k$  where  $v$  is Archimedean and all finite places  $v$  at which  $\phi_v(\cdot, s)$  is not  $K_v$  invariant.

Then the normalized Eisenstein Series  $b_{n,S}(s)E_n^n(g, s, \phi)$  has at most simple poles and these may occur only in the set  $X_n = \{-\varrho_n, 1 - \varrho_n, \dots, \hat{0}, \dots, \varrho_n - 1, \varrho_n\}$  where  $\varrho_n = \frac{1}{2}(n+1)$  ( $\frac{1}{2}(n-1)$  resp.) if  $\varepsilon = -1$  ( $\varepsilon = +1$  resp.). Here  $\hat{0}$  means 0 is omitted in the case when  $n$  is odd.

**Corollary.** We let  $\Pi$  be an automorphic irreducible cuspidal representation of  $H_1(\mathbb{A})$  defined in §1. We let  $G = G_0$  be the doubled group and  $S(\Pi) = S$  as given in §0.

The poles of  $L_S(s, \Pi, r)$  are at most simple and these may occur only in the set  $1/2 + X_\ell(X_\ell$  defined relative to the doubled group  $G = G_0$  with  $\ell = \dim V$ ).

Thus to determine whether the poles of  $b_{n,S}(s)E_n^n(g, s, \phi)$  actually occur we must study what the residue representation is. That is we must determine the subspace

$$\operatorname{Res}_{s=s_0} b_{n,S}(s)E_n^n(g, s, \phi) = (E_n^n)^*(g, s_0, \phi)$$

(with  $s_0 \in X_n$ ) as  $\phi$  varies.

We now consider the group  $\operatorname{Sp}_n = G_0$ . Much of the discussion remains valid for  $O(n, n)$  but since we have not checked all the details we restrict  $G_0$  just to  $\operatorname{Sp}_n$ . We consider the oscillator representation  $\omega_{Q,\psi}$  (defined in §0) of  $\operatorname{Sp}_n \times O(Q)$  on the space  $S(M_{mn}(k))(m = \dim Q)$ . In the case where  $k_v$  is non-Archimedean, we determine the space of  $O(Q)$  coinvariants in  $S(M_{mn}(k))$ , which is given by  $S(M_{mn}(k))/\{\text{the span of } \varphi - \omega_{Q,\psi}((\cdot, y))\varphi \text{ as } \varphi \text{ varies in } S(M_{mn}(k))\} \cong S(M_{mn}(k))_{O(Q)}$ . Then we know that the functional  $\varphi \rightsquigarrow \varphi(0)$  has the property that  $f_\varphi(pg) = \langle \delta | \omega_{Q,\psi}(pg, y) \varphi \rangle = |\det m'|^{\frac{\dim Q}{2}} \chi_Q(\det m') \langle \delta | \varphi \rangle$  (where  $p = m' \cdot v$  with  $m' \in GL_n$ ). Thus  $f_\varphi \in I_v(\frac{m}{2} - \frac{n+1}{2}, \chi_Q)$ . Here  $I_v(s, \chi_Q)$  denotes the local representation of  $\operatorname{Sp}_n$  given by  $\operatorname{Ind}_{P_n}^{S_p}(|\det g|^s \otimes \chi_Q \otimes 1_{U_n})$  (normalized induction), where  $\chi_Q$  is the quadratic character given by the Hilbert symbol  $\langle |\Delta_Q| \rangle$  with  $\Delta_Q$  = discriminant of  $Q$ . In fact the structure of  $S(M_{mn}(k))_{O(Q)}$  as a  $\operatorname{Sp}_n$  module is governed by the space of  $f_\varphi$ .

**Proposition 3.1** [R-I, K-R-III]. Every  $O(Q)$  invariant functional on the space  $S(M_{mn}(k))$  factors through the map  $S(M_{mn}(k)) \rightsquigarrow I_v(\frac{m}{2} - \frac{n+1}{2}, \chi_Q)$  given by  $\varphi \rightsquigarrow f_\varphi$ . In the case  $k = k_v$  is non-Archimedean then  $S(M_{mn}(k))_{O(Q)}$  is equivalent as an  $\operatorname{Sp}_n$  module to  $\{\operatorname{Span} f_\varphi | \varphi \in S(M_{mn}(k))\} = R^n(Q)$ . In case  $k_v = \mathbb{R}$  we let  $R^n(Q) = \{\operatorname{Span} f_\varphi | \varphi \in S(M_{mn}(\mathbb{R})) \text{ and } \varphi \text{ is } U(n) \text{ finite}\}$ ,  $U(n) = \text{maximal compact subgroup of } \operatorname{Sp}_n(\mathbb{R})$ . Moreover for each place  $v$ ,  $R^n(Q)$  is irreducible if  $\dim Q \leq n+1$ .

**Remark.** 3.1. This is the local version of the Siegel-Weil formula.

Then given a global quadratic form  $Q$  over  $k$  we define  $R^n(Q) = \otimes_v R^n(Q_v)$ , which is an admissible representation of  $\operatorname{Sp}_n(\mathbb{A})$ .

The representation  $R^n(Q)$  has the following automorphic structure.

**Theorem 3.2** [K-R-VI].

Let  $\dim Q < n + 1$ . Then

(i)  $\dim \text{Hom}_{\text{Sp}_n(\mathbb{A})}(R^n(Q), \mathcal{A}(\text{Sp}_n(\mathbb{A}))) \leq 1$ ; In fact

(ii) there exists a nonzero embedding  $r_Q^n$  of  $R^n(Q)$  into  $\mathcal{A}(\text{Sp}_n(\mathbb{A}))$ . Moreover the image  $r_Q^n(R^n(Q)) \subseteq L^2_{\text{res}}(\text{Sp}_n(\mathbb{A}))$ , the residual spectrum of  $L^2(\text{Sp}_n(\mathbb{A}))$ , except in the case when  $Q$  is the split 2 dimensional form.

At this point we say two global quadratic forms  $Q$  and  $Q'$  are complementary if  $\dim Q + \dim Q' = 2n + 2$  and if  $Q \cong A \oplus H_{2r}$ , where  $A$  is anisotropic and  $H_{2r}$  is a direct sum of  $r$  hyperbolic planes, then  $Q' \cong A \oplus H_{2r'}$ . We note that  $Q$  uniquely determines  $Q'$  and vice-versa.

We note that if  $\dim Q \geq n + 1$  and  $Q$  admits a complementary form  $Q'$  then for each place  $v$ ,  $R^n(Q_v)$  admits a unique quotient representation  $R^n(Q'_v)$  (this is basically just a simple case of the Howe duality conjecture for dual reductive pairs).

Thus we can describe the residue representation for  $s_0 \in X_n$  and  $s_0 > 0$ . Indeed we note that  $b_{n,S}(s)$  is analytic and non vanishing at such  $s_0$ . Thus we unambiguously have that

$$\begin{aligned} \text{Subspace spanned by } (E_n^n)^*(g, s_0, \phi) = \\ \text{Subspace spanned by } \underset{s=s_0}{\text{Res}} E_n^n(g, s, \phi) = R_n(s_0) \end{aligned}$$

This makes it possible to give a simple characterization of  $R_n(s_0)$  in the following terms.

**Theorem 3.3** [K-R-VI]. Let  $s = s_0 \in X_n$  and  $s_0 > 0$ . Then as a  $\text{Sp}_n(\mathbb{A})$  module  $R_n(s_0)$  is isomorphic to the direct sum  $\oplus R^n(Q')$  where the direct sum ranges over all classes of quadratic forms  $Q'$  with  $\dim Q' = (n + 1) - 2s_0$  and  $\chi_Q = 1$ .

**Remark.** 3.2. We note that in case  $s_0 = \frac{n+1}{2}$ ,  $R_n(\frac{n+1}{2})$  is the identity representation if  $\chi = 1$  and  $\{0\}$  otherwise. In case  $s_0 = \frac{n-1}{2}$ , then  $R_n(\frac{n-1}{2}) = R^n(H_2)$  if  $\chi = 1$ . Otherwise in the remaining cases  $R_n(s_0)$  is a direct sum of infinitely many  $R^n( )$ .

Using Theorem 3.3 it is possible to show consequences concerning the existence of the poles of standard  $L$  function of the group  $\text{Sp}_n$ .

**Theorem 3.4.** Suppose  $L_S(\Pi, r)$  admits a pole at  $s = \frac{1}{2} + s_0$  where  $s_0 \in X_{2n}$  and  $s_0 > 0$ . Then there exists a form  $Q'$  where  $\dim Q' = 2n + 1 - 2s_0$  and  $\chi_{Q'} = 1$  with the property that relative to the dual pair  $\text{Sp}_n \times \text{O}(Q')$  we can find  $\varphi \in S(M_{\dim Q', n}(\mathbb{A}))$  and a function  $f_\Pi \in \Pi$  so that the  $\theta$  integral  $\theta_\varphi(f_\Pi) \neq 0$ .

**Remark.** 3.3. We show in §4 that the actual set of possible poles of  $L_S(s, \Pi, r)$  for  $\text{Re}(s) > 0$  is the set  $s \in \{1, 2, \dots, [\frac{n}{2}] + 1\}$ .

**Remark.** 3.4. If  $L_S(s, \Pi, r)$  admits a pole at  $s = \frac{1}{2} + s_0$ , then there may be more than one  $Q'$  so that  $\theta_\varphi(f_\Pi) \neq 0$ .

## §4. Siegel-Weil Formula

We consider the global version of Proposition 3.1. This constitutes a generalized Siegel-Weil formula, describing the residue of Eisenstein series in terms of  $\theta$  series. First we recall the well known version of this formula established in [W] and extended in [K-R-I] and [K-R-II].

We define a particular type of global section built from  $f_\varphi$  for  $\varphi \in S(M_{mn}(k))$ . We consider  $f_\varphi(g) = \prod_v f_{\varphi_v}(g_v)$  where  $\varphi = \otimes \varphi_v$ . Then  $f_\varphi \in I(\frac{m}{2} - \frac{n+1}{2}, \chi_Q)$ . Then we consider the character on  $P_n$  given by  $p = m' \cdot u \rightsquigarrow |\det m'|^{s-s_0}$  ( $s_0 = \frac{m}{2} - \frac{n+1}{2}$ ). We define the function on  $\mathrm{Sp}_n(\mathbb{A})$

$$e^{sH_P(g)} = |\det m'|^{s-s_0}$$

where  $g = m'uk$  relative to the Iwasawa decomposition. Then we form

$$F_\varphi^n(g, s) = f_\varphi(g) e^{sH_P(g)}$$

Now  $F_\varphi^n$  is a standard  $\phi(, s)$  section defined in §1. Hence we may form Eisenstein series having certain remarkable properties as given by the Siegel-Weil formula.

**Theorem 4.1** [W, K-R-I, K-R-II].

Let either  $Q$  be anisotropic or  $\dim Q - \text{Witt index } (Q) > n+1$ . Then the integral

$$\int_{O(Q)(k) \backslash O(Q)(\mathbb{A})} \theta_\varphi(x, y) dy$$

is absolutely convergent. The Eisenstein series  $E_n^n(g, s, F_\varphi^n)$  is holomorphic at  $s = s_0 = \frac{m}{2} - \frac{n+1}{2}$  and defines an intertwining map of  $R^n(Q)$  into  $\mathcal{A}(\mathrm{Sp}_n(\mathbb{A}))$ . Moreover we have the identity

$$E_n^n(x, s_0, F_\varphi^n) = c_Q \int_{O(Q)(k) \backslash O(Q)(\mathbb{A})} \theta_\varphi(x, y) dy$$

with  $c_Q \neq 0$  depending only on  $Q$ .

The point in extending the Siegel-Weil formula is to have a procedure of regularizing the  $\theta$  integral above.

At this point we look at the dual pair  $\mathrm{Sp}_n \times O(Q)$  with  $\dim Q \leq 2n$  and  $\chi_Q = 1$ .

Let  $v = \infty$  be an Archimedean real place. Then the local form  $Q_\infty$  is equivalent to  $H_{2r} \oplus V_0$  where  $V_0$  is anisotropic and  $H_{2r}$  is a direct sum of  $2r$  hyperbolic planes. We consider the family of oscillator representations  $\omega_{H_{2i} \oplus V_0, \psi}$  (with  $i \leq r$ ).

We let  $\mathcal{Z}_{\mathrm{Sp}_n} = \mathcal{Z}$  be the center of the universal enveloping algebra of  $\mathrm{Sp}_n(\mathbb{R})$ .

We define the support ideal  $I_{Q_\infty}$  of  $\omega_{Q_\infty}$  as

$$\{\xi \in \mathcal{Z} | \omega_{H_{2i} \oplus V_0, \psi}(\xi)(\varphi_\infty) = 0 \text{ for all } \varphi_\infty \in S(M_{\dim(H_{2i} \oplus V_0), n}(\mathbb{R})) \text{ for } i = r-1\}.$$

Then it is straightforward to verify that  $I_{Q_\infty} \neq \mathcal{Z}$  (since  $\dim Q \leq 2n$ ). Moreover it is also easy to check that the space of  $\theta$  kernels on the group  $\mathrm{Sp}_n \times O(Q)$

$$\{\theta_{\omega_{Q_\infty, \psi_\infty}(\xi_\infty)(\varphi)}(x, y) | \xi_\infty \in I_{Q_\infty}\}$$

is rapidly decreasing in the  $O(Q)$  variable ( $\dim Q \leq 2n$  and  $Q \neq H_2$  split form).

Then as a substitute for the  $\theta$  integral we use the regularized  $\theta$  integral with an element  $\xi_\infty \in I_{Q_\infty}$ . That is,

$$i_{\xi, Q}(\varphi) = \int_{O(Q)(k) \backslash O(Q)(\mathbb{A})} \theta_{\omega_{Q_\infty, \psi_\infty}(\xi_\infty)(\varphi)}(x, y) dy.$$

Then  $i_{\xi, Q} \in \text{Hom}_{\text{Sp}_n(\mathbb{A})}(R^n(Q), \mathcal{A}(\text{Sp}_n(\mathbb{A}))$ . But we also have

**Lemma 4.1.** *If  $\dim Q - \text{Witt index } (Q) \leq n + 1$  and  $2n \geq \dim Q > n + 1$  then  $i_{\xi, Q}$  factors through the space  $R^n(Q')$ ,  $Q' =$  the complementary form to  $Q$ .*

Moreover if  $\dim Q \leq n + 1$  and  $Q \neq H_2$ , then the map  $s \rightsquigarrow E_n^n(g, s, F_\varphi^n)$  is analytic at  $s = s_0 = \frac{m}{2} - \frac{n+1}{2}$ . Also  $\varphi \rightsquigarrow \underset{s=s_0}{\text{val}} E_n^n(g, s, F_\varphi^n)$  defines an element in  $\text{Hom}_{\text{Sp}_n(\mathbb{A})}(R(Q), \mathcal{A}(\text{Sp}_n(\mathbb{A}))$ .

If  $2n \geq \dim Q > n + 1$  and  $\dim Q - \text{Witt index } (Q) \leq n + 1$  then  $\varphi \rightsquigarrow \underset{s=s_0}{\text{Res}} E_n^n(g, s, F_\varphi^n)$  (with  $s_0 = \frac{m}{2} - \frac{n+1}{2}$ ) defines a intertwining map which factors through  $R^n(Q')$ . This determines a nonzero element in  $\text{Hom}_{\text{Sp}_n(\mathbb{A})}(R^n(Q'), \mathcal{A}(\text{Sp}_n(\mathbb{A}))$ .

By the Uniqueness Principle of Theorem 3.2 we have the following Siegel-Weil type identity.

**Theorem 4.2.** *Let  $s_0 = \frac{\dim Q}{2} - \frac{n+1}{2}$ .*

(1) *Let  $\dim Q \leq n + 1$ , with  $Q \neq H_2$ . Then there is a constant  $c_\xi$  so that*

$$c_\xi E_n^n(g, s_0, F_\varphi^n) = i_{\xi, Q}(\omega_{Q, \psi}(g)\varphi)$$

(2) *Let  $2n \geq \dim Q > n + 1$  and  $\dim Q - \text{Witt index } (Q) \leq n + 1$ . Then there is a constant  $c'_\xi$  so that*

$$c'_\xi \underset{s=s_0}{\text{res}} E_n^n(g, s, F_\varphi^n) = i_{\xi, Q}(\omega_{Q, \psi}(g)\varphi)$$

The goal then is to use Theorems 4.1 and 4.2 to glean information about special values of the  $L$  functions discussed above.

We indicate one type of result in this direction.

We fix a “Witt tower” of quadratic forms  $Q_0, Q_0 \oplus H_2 = Q_2, \dots, Q_0 \oplus H_{2r} = Q_{2r}$  where  $Q_0$  is anisotropic (even dimensional) and  $H_{2r}$  is a direct sum of  $r$  hyperbolic planes.

We let  $X_{Q_{2i}} = \{f \in L^2_{\text{cusp}}(\text{Sp}_n(\mathbb{A})) | \theta_\varphi(f) = 0 \text{ for all } \varphi \in S(M_{\dim Q_{2i}, n}(\mathbb{A}))\}$ . Then  $X_{Q_{2r}}$  is a  $\text{Sp}_n(\mathbb{A})$  module and we define inductively the spaces  $Y_{Q_{2i}}$  = the perpendicular complement of the space  $X_{Q_{2i}} \cap X_{Q_{2i-2}} \cap \dots \cap X_{Q_0}$  in  $X_{Q_{2i-2}} \cap \dots \cap X_{Q_0}$ . Here perpendicular complement is taken relative to the Hermitian pairing given by the Petersson inner product on  $L^2_{\text{cusp}}(\text{Sp}_n(\mathbb{A}))$ . Then we have an orthogonal direct sum decomposition of  $\text{Sp}_n(\mathbb{A})$  modules of  $L^2_{\text{cusp}}(\text{Sp}_n(\mathbb{A})) = Y_{Q_0} \oplus Y_{Q_2} \oplus$

$\cdots \oplus Y_{Q_{4n}}$ . The space  $Y_{Q_{2i}} = \{0\}$  if  $\dim Q_{2i} < n$  ([PS-R-IV]). Also, if  $f_\Pi \in \Pi$ , an irreducible automorphic cuspidal representation which occurs in  $Y_{Q_n}$ , then there exists  $\varphi \in S(M_{mn}(\mathbb{A}))$  ( $m = \dim Q_{2i}$ ) ([R-I]) so that  $\theta_\varphi(f_\Pi) \neq 0$ . In fact,  $\theta_\varphi(f_\Pi)$  is a cusp form on  $O(Q_{2i})(\mathbb{A})$ .

On the other hand if  $f \in Y_{Q_{2i}} \oplus \cdots \oplus Y_{Q_{4n}}$  then whether  $\theta_\varphi(f)$  (for  $\varphi \in S(M_{m,n}(\mathbb{A}))$ ) is zero or not is measured by the inner product formula developed in [R-II] and [R-III]. Indeed, using Theorems 4.1 and 4.2 we show the following identity:

**Proposition 4.1.** *Let  $f_\Pi \in \Pi$  be chosen as in the beginning of §2 where  $\Pi \subset Y_{Q_{2i}} \oplus \cdots \oplus Y_{Q_{4n}}$ .*

(1) *Then if either  $\dim Q_{2i} \leq 2n + 1$  or  $\dim Q_{2i} - i > 2n + 1$*

$$\|\theta_{\varphi_1}(f_\Pi)\|_{O(Q_{2i})(\mathbb{A})}^2 = c'_{Q_{2i}} \underset{s=s_0}{\text{val}} Z_{2n}(f_{\Pi_1} \otimes \bar{f}_{\Pi_1}, F_{\varphi_1}^{2n}(\cdot, s))$$

where  $\varphi_1^*$  is a canonically determined function (coming from  $\varphi_1$ ) in  $S(M_{\dim Q_{2i}, 2n}(\mathbb{A}))$  (the doubled Weil representation space of  $\text{Sp}_{2n} \times O(Q_{2i})$ ). Here  $s_0 = \frac{\dim Q_0}{2} + i - (n + \frac{1}{2})$ . Here  $\|\cdot\|_{O(Q_{2i})(\mathbb{A})}^2$  represents the Petersson norm of a function relative to a suitably normalized measure on  $O(Q_{2i}(k) \backslash O(Q_{2i})(\mathbb{A}))$ . Also  $c'_{Q_{2i}}$  is a nonzero constant independent of  $\varphi$  and  $f_\Pi$ .

(2) *The inner product formula of (1) is also valid in the case where  $\dim Q_{2i} - i \leq 2n + 1$  and  $2n + 1 < \dim Q_{2i} < 4n$ , but now  $\text{val}$  is replaced by  $\underset{s=s_0}{\text{Res}}$ .*

Thus, using Proposition 2.1, it is possible to relate the Petersson norm of  $\theta_{\varphi_1}(f_\Pi)$  to the special value of the  $L_s$  function defined in §2. The main problem that remains in getting an exact relation is to have control on the bad local factors, i.e. where the data is not spherical in  $Z_v(\cdot \cdot \cdot)$ . The point here is to be able to determine the local factor  $Z_v(\cdot \cdot \cdot)$  when the local data has the form  $f_\varphi$  as given above.

However in any case using Proposition 2.1 and a much simpler version of the identity in Proposition 4.1 we deduce the statement about the poles of the  $L$  functions given in Theorem 3.4. Moreover Remark 3.3 follows from the fact that  $Y_{Q_{2i}} = 0$  for  $\dim Q_{2i} < n$ .

## References

- [B-J] Borel, A., Jacquet, H.: Automorphic forms and automorphic representations. Proc. Symp. Pure Math. 33, vol. 1, 1979
- [G-S] Gelbart, S., Shahidi, F.: Analytic properties of automorphic  $L$  functions. Perspectives in Math., Academic Press, 1988
- [Kim] Kimelfeld, B.: Homogeneous domains on flag manifolds. J. Math. Anal. Appl. **121** (1987) 506–588
- [K-R-I] Kudla, S., Rallis, S.: On the Weil-Siegel formula. J. Reine Angew. Math. **387** (1988) 1–68
- [K-R-II] Kudla, S., Rallis, S.: On the Weil-Siegel formula II. J. Reine Angew. Math. **391** (1988) 65–84

- [K-R-III] Kudla, S., Rallis, S.: Degenerate principal series and invariant distributions. *Israel J. Math.* **69** (1990) 25–45
- [K-R-IV] Kudla, S., Rallis, S.: Poles of Eisenstein series and  $L$  functions. Preprint 1989
- [K-R-V] Kudla, S., Rallis, S.: Ramified degenerate principal series. In preparation, 1990
- [K-R-VI] Kudla, S., Rallis, S.: A regularized Siegel-Weil formula: the first term identity. In preparation, 1990
- [PS-R-I] Piatetski-Shapiro, I., Rallis, S.:  $\varepsilon$  factor of representations of classical groups. *Proc. Natl. Acad. Sci.* **83** (1986) 4589–4593
- [PS-R-II] Piatetski-Shapiro, I., Rallis, S.:  $L$  functions for classical groups. (Lecture Notes in Mathematics, vol. 1254). Springer, Berlin Heidelberg New York 1987
- [PS-R-III] Piatetski-Shapiro, I., Rallis, S.: Rankin Triple  $L$  functions. *Comp. Math.* **64** (1987) 31–115
- [PS-R-IV] Piatetski-Shapiro, I., Rallis, S.: A new way to get Euler products. *J. Reine Angew. Math.* **392** (1988) 110–124
- [PS-R-S] Piatetski-Shapiro, I., Rallis, S., Schiffmann, G.: Rankin-Selberg integrals for  $G_2$ . To appear in *Amer. J. Math.*
- [Rab] Rabau, P.: Preprint 1990
- [R-I] Rallis, S.: On the Howe duality conjecture. *Comp. Math.* **51** (1984) 333–399
- [R-II] Rallis, S.: Injectivity Properties of liftings associated to Weil representation. *Comp. Math.* **52** (1984) 139–169
- [R-III] Rallis, S.:  $L$  functions and the oscillator representation. (Lecture Notes in Mathematics, vol. 1245). Springer, Berlin Heidelberg New York 1987
- [W] Weil, A.: Sur la formule de Siegel dans la theorie des groupes classiques. *Acta. Math.* **113** (1965) 1–87



# Iteration of Polynomial Automorphisms of $\mathbf{C}^2$

Eric Bedford

Department of Mathematics, Indiana University, Swain Hall East  
Bloomington, IN 47405, USA

## § 1. Introduction

We will consider mappings  $f = (f_1, f_2) : \mathbf{C}^2 \rightarrow \mathbf{C}^2$  such that  $f_1$  and  $f_2$  are holomorphic polynomials. A polynomial automorphism is a polynomial mapping which is 1-to-1 and onto; it follows that the inverse  $f^{-1}$  is also polynomial. We will discuss such mappings from the point of view of dynamical systems. That is, we will be primarily concerned with the forms of limiting and/or recurrent behavior that the iterates  $\{f, f^2 = f \circ f, \dots, f^{on}, \dots\}$  of such mappings can exhibit. Our purpose here is to present some work that we have done on these problems in collaboration with J. Smillie and which has been written up in [BS1, 2, 3]. Lack of time prevents us from discussing more recent work on entropy, Lyapunov exponents, and ergodicity (see [BS4]).

Analogous problems have been studied for polynomial mappings  $p : \mathbf{C} \rightarrow \mathbf{C}$  in one complex variable. The starting point is the Julia set  $J_p$ , which is the complement of the set where the iterates  $\{p^n : n = 1, 2, 3, \dots\}$  are well behaved, i.e.  $\mathbf{C} - J_p$  is the complement of the largest open set where  $\{p^n\}$  forms a normal family.  $J_p$  turns out to be the boundary of  $K_p^+ = \{z : p^n(z) \text{ is bounded for } n = 1, 2, 3, \dots\}$  is a nonempty, compact, invariant set which carries all of the interesting dynamics of  $p$ . The fundamental properties of  $p$  on  $J_p$  were developed in the classical work of Fatou and Julia. One of the basic results of the theory is:

$$\text{The repelling periodic points for } p \text{ are dense in } J_p. \quad (1)$$

Here we will consider the potential theoretic approach to the dynamical properties of a polynomial  $p$  (cf. [Br, T]). This approach is well adapted to obtaining results of a more precise quantitative nature. The equilibrium measure  $\mu$  of  $J_p$ , which is fundamental in potential theory, turns out to coincide with the unique measure of maximal entropy. Two results of Brolin [Br] show that  $\mu$  arises as the limit of two sequences of sets of points. First, it is the limit of averages of point masses on the preimages of a point, i.e.

$$\mu = \lim_{n \rightarrow \infty} (\deg)^{-n} \sum_{a \in \{p^{-n}(z_0)\}} \delta_a, \quad (2)$$

where the limit is taken with respect to the weak topology of measures. Second, it is the limit of the average of point masses of periodic points, i.e.

$$\mu = \lim_{n \rightarrow \infty} (\# \text{Per}(n))^{-1} \sum_{a \in \text{Per}(n)} \delta_a, \quad (3)$$

where  $\text{Per}(n) = \{z \in \mathbf{C} : p^n(z) = z\}$  is the set of periodic points whose periods divide  $n$ .

While the potential theoretic approach gives useful and powerful tools, it applies only to polynomial mappings. Although much of the theory of Fatou and Julia applies also to rational mappings, potential-theoretic methods do not apply in this case. The equilibrium measure of  $J_p$  is not equal to the measure of maximal entropy unless  $p$  is a polynomial (see [Lo]). Although problems of iteration have been studied for some classes of entire functions, the theory is quite different; for instance,  $J_p$  is not the boundary of  $K_p^+$ .

A prototypical example of the polynomial automorphisms which we will consider is the mapping  $f : \mathbf{C}^2 \rightarrow \mathbf{C}^2$  given by

$$f(x, y) = (y, p(y) - ax), \quad (4)$$

where  $p(y)$  is a one-variable polynomial of degree  $d > 1$ . In fact, as will be shown in §2, there is no essential loss of generality in considering only these maps.

There are several sources of motivation for studying the iteration of automorphisms of  $\mathbf{C}^2$ . One is to generalize to the 2-dimensional case the results that have been obtained for complex polynomial maps. The two basic properties of a 1-variable map: that it is  $d$ -to-1 ( $d > 1$ ) and that it is conformal, do not hold for automorphisms of  $\mathbf{C}^2$ . However, the fact that a 1-dimensional map  $p : \mathbf{C} \rightarrow \mathbf{C}$  is a proper mapping of degree  $d$  generalizes to the observation that the mapping on homology given by (10) below is multiplication by  $d$ . Conformality of  $p$  is partially replaced by the fact that  $f$  is conformal when restricted to leaves of the stable/unstable foliations.

We define the *stable set* of a point  $p$  as

$$W^s(p, f) = \{z \in \mathbf{C}^2 : \lim_{n \rightarrow \infty} \text{dist}(f^n(z), f^n(p)) = 0\}. \quad (5)$$

We say that  $p$  is a *periodic point* for  $f$  if  $f^m(p) = p$  for some  $m \geq 1$ , and we call the smallest such  $m$  the *period* of  $p$ . For a periodic point  $p$ , we let  $\lambda_1, \lambda_2$  denote the eigenvalues of  $Df^m(p)$ . If  $|\lambda_1|, |\lambda_2| < 1$ , then  $p$  is a *sink*, and it is well known that  $W^s(p, f)$  is an open set containing  $p$ , which is called the *basin of attraction* of  $p$ . A proper open subset  $\Omega \subset \mathbf{C}^2$  which is biholomorphically equivalent to  $\mathbf{C}^2$  is called a *Fatou-Bieberbach domain*. It is classically known that if  $p$  is a sink, then the basin  $B = W^s(p, f)$  is biholomorphically equivalent to  $\mathbf{C}^2$ , and in many cases, e.g. if  $f$  has more than one sink,  $B$  is a Fatou-Bieberbach domain which is not dense in  $\mathbf{C}^2$ . Basins of attraction were the original examples of Fatou-Bieberbach domains. The geometry of Fatou-Bieberbach domains in general is intriguing but not well understood, and the ones that arise as basins of attraction seem to be the most approachable. For instance, it was shown in [BS2] that a polynomial basin  $B = W^s(p, f)$  cannot be “small” enough to avoid an algebraic variety, or “large” enough to contain one. In other words, if  $V$  is a 1-dimensional algebraic variety, then

$$B \cap V \neq \emptyset \quad \text{and} \quad V \neq \overline{B}.$$

Another motivation comes from polynomial automorphisms of  $\mathbf{R}^2$ . Hénon [Hn1, 2] showed that the automorphism  $g : \mathbf{R}^2 \rightarrow \mathbf{R}^2$  given by

$$g(x, y) = (y, y^2 + c - ax) \quad (6)$$

can possess complicated dynamics. The stable set of a subset  $A \subset \mathbf{C}^2$  is given by

$$W^s(A, g) = \{p \in \mathbf{C}^2 : \lim_{n \rightarrow \infty} \text{dist}(g^n(p), A) = 0\}. \quad (7)$$

Let us call a set  $A$  an *attractor* for  $g$  if  $A$  is compact and invariant, and  $W^s(A)$  contains a neighborhood of  $A$ . It was shown recently (see [BC]) that  $g$  has a “strange attractor” for certain values of the parameters  $a$  and  $c$ . “Strange” implies that the attractor  $A$  is not a union of sink orbits. One difference between the real and complex cases is worth pointing out here. Although this set  $A$  is an attractor for  $g$  in  $\mathbf{R}^2$ , it is not an attractor, when considered in  $\mathbf{C}^2$ . In fact a normal families argument shows that any attractor for  $g$  in  $\mathbf{C}^2$  is a union of sink orbits. Although it is not clear exactly what relationship holds between the mapping  $g$  and the corresponding mapping extended to  $\mathbf{C}^2$ , it is felt that the study of the map in  $\mathbf{C}^2$  should be analogous to approaching the study of polynomial maps  $p : \mathbf{R} \rightarrow \mathbf{R}$  by the study of the same polynomial in the complex domain. When the theory is pursued in the complex domain, it seems to be more complete and show a more stable or continuous dependence on parameters.

## § 2. Elementary Mappings

The map  $f$  and a conjugate  $h^{-1} \circ f \circ h$  will have the same properties under iteration. Friedland and Milnor [FM] have shown that the set of polynomial automorphisms naturally divides into two sets of equivalence classes under conjugation by polynomial automorphisms. One of these classes is called the “elementary” automorphisms, and the other consists of finite compositions of the form

$$f = f_1 \circ \dots \circ f_m, \quad (8)$$

where

$$f_j(x, y) = (y, p_j(y) - a_j x) \quad (9)$$

and  $p_j(y) = y^{d_j} + c_{d_j-2}y^{d_j-2} + \dots$  is a monic polynomial of degree  $d_j > 1$  in which the  $y^{d_j-1}$  coefficient vanishes.

The dynamics of the elementary maps has been studied in detail by Friedland and Milnor. They have shown that an elementary map  $f$  has rather simple dynamics;  $f$  has periodic points of only finitely many periods, the nonwandering set is quite simple, and  $f$  is conjugate to an isometry on the nonwandering set.

Thus the rest of this talk will be concerned with mappings of the form (8). If we define the sets

$$V^\pm = \{(x, y) \in \mathbf{C}^2 : |x| \gtrless |y|, \max|x|, |y| > \kappa\}$$

$$V = \{\max|x|, |y| < \kappa\}$$

then for  $\kappa$  large enough

$$\begin{aligned} f^{-1}(V^+) &\subset V^+ \quad \text{and} \quad f^{-1}(V^+ \cup V) \subset V^+ \cup V \\ f(V^-) &\subset V^- \quad \text{and} \quad f(V^- \cup V) \subset V^- \cup V. \end{aligned}$$

It is evident, then that for a point  $p$  in  $V^+$  (or  $V^-$ )  $f^n(p)$  tends to infinity as  $n$  tends to  $-\infty$  (or  $+\infty$ ). Thus all recurrent behavior of  $f$  takes place in the bounded set  $V$ . (This is useful, since there does not seem to be a natural extension of  $f$  to a compactification of  $\mathbf{C}^2$ ; the one point compactification of  $\mathbf{C}^2$ , for instance, cannot be given a complex structure.)

Although  $f$  is an automorphism, the behavior of  $p$  in the large is reflected by the fact that the mapping on homology

$$f_* : H_2(V^- \cup V, V^-) \rightarrow H_2(V^- \cup V, V^-) \quad . \quad (10)$$

corresponds to multiplication by  $d$ , i.e.  $f_*(\tau) = (d)\tau$ .

### § 3. The Sets $K^\pm$

The approach adopted by Hubbard and Oberste-Vorth [HO] is to consider the sets

$$K^\pm = \{q \in \mathbf{C}^2 : f^{\pm n}(q) \text{ is bounded for } n = 1, 2, 3, \dots\}$$

where the iterates stay bounded in forward/backward time. Also of interest are the sets

$$J^\pm = \partial K^\pm$$

$$K = K^+ \cap K^- \quad \text{and} \quad J = J^+ \cap J^-.$$

It is evident that with  $V^\pm$ ,  $V$  as in the previous section, then  $K^\pm \subset V \cup V^\pm$ .

A sequence  $\{q_n\}$  is said to be an  $\varepsilon$ -orbit if  $\text{dist}(q_{n+1}, f(q_n)) < \varepsilon$  holds for all  $n = 1, 2, 3, \dots$ . A point  $q$  is *chain recurrent* if for any  $\varepsilon > 0$  there is an  $\varepsilon$ -orbit  $\{q_n\}$  with  $q = q_1 = q_{1+N}$  for some  $N$  and all  $j = 1, 2, 3, \dots$ . The set of chain recurrent points is denoted by  $R(f)$ . Since the iterates of any point outside of  $K$  tend to infinity in either positive or negative time, we see that  $R(f) \subset K$ .

In contrast to the case of mappings of one variable, where the Julia set can be either connected or a Cantor set, *the sets  $K^\pm$  and  $J^\pm$  are always connected* (cf. [BS1], Theorem 7.2.) The set  $J = J^+ \cap J^-$ , on the other hand, may be either connected or disconnected.

Hubbard and Oberste-Vorth [H, HO] have studied  $U^\pm := \mathbf{C}^2 - K^\pm$  in some detail for mappings of the form (6) with  $a$  and  $c$  complex. For all  $a, c$ , the sets  $U^\pm$  are homeomorphic to  $(S^3 - \Sigma) \times \mathbf{R}$ , where  $\Sigma$  is a solenoid, and the fundamental group  $\pi_1(U^\pm) \cong \mathbf{Z}[\frac{1}{2}]$  is not finitely generated. [HO] also makes a detailed study of the topology of  $K^\pm$  for mappings of the form (6) and certain values of  $a, c \in \mathbf{C}$ .

If  $(x, y) \notin K^+$ , then the iterates  $(x_n, y_n) := f^n(x, y)$  become unbounded with asymptotic behavior like  $y_n \sim p(x_n)$ . In the potential theoretic approach, we study  $K^\pm$  and  $\mathbf{C}^2 - K^\pm$  in terms of the rate at which the points  $(x_n, y_n)$  escape to infinity. From the specific polynomial form of  $f$ , we see that for fixed  $x$

$$|f_n(x, y)| = |y|^{d^n} + O(|y|^{d^n-2}) \quad (11)$$

as  $|y| \rightarrow \infty$ . This motivates the definition

$$G^\pm(x, y) = \limsup_{n \rightarrow +\infty} (d^{-n}) \log(|x_{\pm n}| + |y_{\pm n}|).$$

Using (11) we may derive the facts

- $\{G^\pm = 0\} = K^\pm$ .
- $G^\pm$  is pluriharmonic on  $K^\pm$ .
- $G^+(x, y) = \log^+ |y| + o(1)$  for  $x$  fixed and  $y \rightarrow \infty$ .
- $G^-(x, y) = \log^+ |x| - \log |a| + o(1)$  for  $y$  fixed and  $x \rightarrow \infty$ .  
It is immediate from the definition that
- $G^\pm \circ f = (d^{\pm n})G^\pm$ .

Further,  $G^\pm$  is continuous on  $C^2$ ; and by a more subtle argument it may be shown to be Hölder continuous (see [FS]).

From these observations, we see that  $G^\pm$  is the pluriharmonic Green function of  $K^\pm$ . In fact, if  $T$  denotes any complex line (i.e. a 1-dimensional subspace of  $C^2$ ), then the restriction  $G^\pm|_T$  is the classical Green function of  $K^\pm \cap T$  inside the Riemann surface  $T$  (which may be identified with  $C$ ).

## § 4. Potential Theory in $C^2$

Although the Laplacian  $\Delta$  is sometimes useful in problems of several complex variables, it is not an invariant operator. The defect of this lack of invariance shows up in problems of iteration, so we will use  $dd^c$ , which is a biholomorphically invariant version of  $\Delta$ . Similarly, the usual subharmonic functions, which are defined by the condition  $\Delta v \geq 0$ , are not preserved under holomorphic transformations. The class of subharmonic functions which remain subharmonic under holomorphic changes of coordinates is the class of *plurisubharmonic* or *psh* functions, and it is this class which will be useful to us.

By  $\mathcal{D}_{(p,q)}$  we will denote the set of test  $(p, q)$ -forms, i.e. the forms which may be written as

$$\alpha = \sum_{1 \leq i_1 < \dots < i_p \leq n} \sum_{1 \leq j_1 < \dots < j_q \leq n} \alpha_{i_1, \dots, i_p, j_1, \dots, j_q} dz_{i_1} \wedge \dots \wedge dz_{i_p} \wedge d\bar{z}_{j_1} \wedge \dots \wedge d\bar{z}_{j_q}$$

where each  $\alpha_{i_1, \dots, i_p, j_1, \dots, j_q}$  is a smooth function with compact support. The dual of this space is the space of  $(p, q)$  currents, which is denoted as  $\mathcal{D}'_{(p,q)}$ . For an example of a  $(1,1)$  current on  $C^2$  let us consider a 1-dimensional closed complex submanifold (or subvariety)  $M$ . The *current of integration*  $[M]$  acts on a  $(1,1)$ -form  $\varphi$  by integration:  $[M](\varphi) = \int_M \varphi$ . A current may be thought of loosely as a  $(p, q)$ -form with distributions as coefficients. In case the current is positive, then the coefficients may be shown to be regular Borel measures. (See [L1] for further details.)

We consider the operator

$$dd^c : \mathcal{D}'_{(0,0)} \rightarrow \mathcal{D}'_{(1,1)},$$

which is given by

$$dd^c u = 2\sqrt{-1} \sum_{i,j=1}^2 \frac{\partial^2 u}{\partial z_i \partial \bar{z}_j} dz_i \wedge d\bar{z}_j.$$

A function  $u$  is psh if it is upper semicontinuous and if  $dd^c u$  is a positive  $(1,1)$  current. The functions  $G^\pm$  defined above are psh, and  $\mu^\pm := dd^c G^\pm$  are positive, closed currents which satisfy

$$f^* \mu^\pm = (d)^{\pm 1} \mu^\pm.$$

We will call  $\mu^\pm$  the *stable/unstable currents*. It it easily seen that the support of  $\mu^\pm$  is exactly  $\partial K^\pm$ .

Since  $\mu^\pm$  are  $(1,1)$  currents, they, unlike measures, do not act directly on sets. But if  $T$  is a locally closed 1-dimensional complex manifold, then we may obtain a measure on  $T$  by restriction:  $\mu^\pm|_T := (dd^c)_T(G^\pm|_T)$ . That is, we let  $G^\pm|_T$  denote the restriction of  $G^\pm$  to  $T$ , and we let  $(dd^c)_T$  denote the induced Laplacian on  $T$ , so  $\mu^\pm|_T$  acts as a measure on subsets of  $T$ .

The first Theorem of Brolin mentioned above concerns the preimages of a point under iteration. Our way of adapting this to our situation is to work with  $(1,1)$  currents. We consider, instead of points, an algebraic variety  $V$  of codimension 1, so that  $[V]$  is a  $(1,1)$  current. The pullback of the current of integration corresponds to the preimage under  $f$  and is given by  $f^{n*}[V] = [f^{-n}V]$ . Brolin's Theorem in this case becomes:

**Theorem [BS1,2].** *If  $V$  is an algebraic variety of codimension 1 in  $\mathbf{C}^2$ , then*

$$\lim_{n \rightarrow \infty} (d^{-n}) f^{n*}[V] = c \mu^+ \quad (12)$$

*holds for some constant  $c > 0$ , with the convergence being taken in the sense of currents.*

(Another version of this Theorem was recently given in [FS].) The proof of this Theorem is based on the Poincaré-Lelong formula for currents

$$[\{x = 0\}] = \frac{1}{2\pi} dd^c \log |x|.$$

Similarly, if  $V = \{h = 0\}$  is an algebraic variety defined by a polynomial  $h$ , then  $dd^c \log |h|$  is a constant (positive) multiple of the current of integration on  $V$ . In vague outline, the proof of the theorem proceeds by showing that

$$\lim_{n \rightarrow \infty} \log |h(f^n(x, y))| = \text{const. } G^+(x, y) \quad (13)$$

holds locally in  $L^1$ . It follows, then, by applying the Poincaré-Lelong identity, that we obtain the convergence of the desired currents.

N. Sibony proposed working directly with invariant measures. If  $S$  is a positive, closed,  $(1,1)$  current on  $\mathbf{C}^2$ , and if  $u$  is a bounded, psh function, then we may form the wedge product  $dd^c u \wedge S$ , which will be a  $(2,2)$  current, and which acts on a test form  $\varphi$  according to the formula

$$(dd^c u \wedge S)(\varphi) := \int u dd^c \varphi \wedge S.$$

This formula is, formally, just an integration by parts, and the integral on the right hand side is justified by representing  $S$  as a form with measure coefficients and performing a (usual) wedge product with the smooth  $(1,1)$  form  $dd^c \varphi$ . In this way, we define

$$\mu := \mu^+ \wedge \mu^-, \quad (14)$$

and it is immediate that

$$f^* \mu = \mu.$$

The two sets of main importance from the dynamical point of view are the chain recurrent set and the nonwandering set. A point  $p$  is *wandering* if there is a neighborhood  $U$  of  $p$  such that  $U \cap f^n U = \emptyset$  for all  $n \neq 0$ , and the *nonwandering set*  $\Omega(f)$  is the set of all nonwandering points. It is generally true that  $\Omega(f) \subset R(f)$ . It is easily seen, too, that the support of an invariant, finite measure is contained in the nonwandering set, i.e.  $\text{spt } \mu \subset \Omega(f)$ .

It has been conjectured by N. Sibony that the analogue of (3) holds for the measure  $\mu$  in (14) and mappings of the form (8), i.e.:

$$\lim_{n \rightarrow \infty} d^{-n} \sum_{z \in \text{Per}(n)} \delta_z = \frac{1}{4\pi^2} \mu. \quad (15)$$

With N. Sibony, we have developed the following formal argument which makes (15) plausible. First  $G := \max\{G^+, G^-\}$  is the psh Green function of the set  $K$ . (See, for instance [Be] for relevant information.) A calculation shows that  $(dd^c G)^2 = \mu$ . If we set

$$u_n = (d^{-n}) \log |f^n(x, y) - f^{-n}(x, y)|,$$

then

$$\lim_{n \rightarrow \infty} u_n = G \quad (16)$$

holds a.e. Further, a standard calculation (cf. [BT1]) shows that

$$(dd^c u_n)^2 = d^{-2n} (2\pi)^2 \sum_{\{a : f^n(a) = f^{-n}(a)\}} \delta_a, \quad (17)$$

which may be interpreted as

$$(dd^c u_n)^2 = (\# \text{Per}(2n))^{-1} (2\pi)^2 \sum_{\{a \in \text{Per}(2n)\}} \delta_a. \quad (18)$$

This information, however, is not sufficient to obtain (15), because (16) does not in general imply that the measures  $(dd^c u_n)^2$  converge to  $(dd^c G)^2$ . (See [C1,2] and [L2] for examples of bad behavior of  $(dd^c)^2$ .) For this approach to work, the limit (16) needs to be taken in a stronger sense.

## § 5. Iteration of Disks

We let  $M$  denote a 1-dimensional, locally closed, complex submanifold of  $\mathbf{C}^2$ , and we study with the current of integration  $[M]$  under pull-backs by the iterates of  $f$ . In order to remove any technical problems that might be introduced by the boundary of  $M$ , we let  $\chi$  be a test function such that  $\text{spt } \chi \cap M$  is compact, and we work with  $\chi[M]$ . The following was obtained in [BS3, Theorem 3].

**Theorem.** *Let  $M$  be as above, and suppose that one of the following holds:*

(i)  $M$  is an open subset of an algebraic curve  $X$ .

(ii)  $M \subset J^+$ .

*If we set  $c = \int_M \chi \mu^-|_M$ , then*

$$\lim_{n \rightarrow \infty} (d^{-n}) f^{*n}(\chi[M]) = \frac{c}{4\pi^2} \mu^+.$$

This result has two interesting consequences for the study of  $K^+$ . The first is

**Theorem.** *Let  $\Omega$  be a connected component of  $\text{int } K^+$ . If  $\Omega \cap J^- \neq \emptyset$ , then  $\partial\Omega = J^+$ .*

In particular, if  $p$  is a sink for  $f$ , then  $p \in J^-$ .

**Corollary.** *If  $p$  is a sink for  $f$ , and if  $B$  is the basin of attraction of  $p$ , then  $\partial B = J^+$ .*

The second application concerns hyperbolic periodic (saddle) points. Let us note that Smillie [Sm] has shown that the map  $f$  has positive topological entropy, so by a result of Katok [K],  $f$  always has saddle points. Such points have stable and unstable manifolds which are complex submanifolds of  $\mathbb{C}^2$ . The following result confirms a conjecture of J.H. Hubbard.

**Theorem.** *If  $p$  is a saddle point for  $f$ , then the stable manifold  $W^s(p)$  is a dense subset of  $J^+$ .*

## § 6. Recurrence

Let us make some remarks on the behavior of  $f$  on  $\text{int } K^+$ . A connected component  $\Omega$  of  $\text{int } K^+$  is said to be a *wandering component* of  $\text{int } K^+$  if  $\Omega \cap f^n(\Omega) = \emptyset$  for  $n \neq 0$ . It is an interesting open question at this time whether there can be a wandering component for  $f$  (cf. the list of questions of Milnor in [Bil]). If there are no wandering components, then all components are periodic. The proof of Sullivan [Su2] that there are no wandering domains in dimension one uses quasiconformal methods, and analogous techniques are not available in the present case.

Here we consider a component  $\Omega$  which is *recurrent*, i.e. there is a point  $z \in \mathbb{C}^2$  such that  $\Omega$  contains a point of the  $\omega$ -limit set  $\omega(z)$ , which is the set of accumulation points of the forward iterates of  $z$ . It is easily shown that if  $\Omega$  is recurrent, then  $f^m(\Omega) = \Omega$  for some  $m \geq 1$ .

Let  $A \subset \mathbb{C}$  denote either the unit disk or an annulus  $\{r_1 < |\zeta| < r_2\}$ . We say that an imbedding  $\varphi : A \rightarrow \mathbb{C}^2$  is *rotation* if there is an irrational number  $a$  such that

$$f(\varphi(\zeta)) = \varphi(e^{i\pi a} \zeta)$$

holds for all  $\zeta \in A$ . In this case we call  $\varphi(A)$  a *rotation domain*.

**Theorem.** *Suppose that  $\Omega$  is a connected component of  $\text{int } K^+$  and that  $f$  decreases volume. If  $\Omega$  is recurrent, then either*

(i)  $\Omega$  is the basin of attraction of a sink.

(ii)  $\Omega$  is the basin of a rotation domain, i.e.  $\Omega = W^s(\varphi(A)) = \bigcup_{\zeta \in A} W^s(\varphi(\zeta))$ .

The possibility of case (ii) was also shown in [FS]. In the second case, it is in fact possible to linearize  $f$  on  $\Omega$ . That is, there is a biholomorphic mapping  $h : A \times \mathbb{C} \rightarrow \Omega$  such that  $h \circ L = f \circ h$ , where  $L(\zeta, w) = (e^{i\pi a} \zeta, \beta w)$ .

We remark that in either case it is evident that  $\Omega \cap J^- \neq \emptyset$ , so by §5 we have  $\partial\Omega = J^+$ .

The remaining possibility is that a component of  $\Omega$  is periodic but not recurrent. Not much is known about this case. T. Ueda [U1,2] has analyzed such domains which occur at a parabolic fixed point.

## § 7. Hyperbolic Mappings

Up to now our discussion has applied equally well to all choices of  $f$ . Now we see what extra information we can derive in the special situation of hyperbolicity. This should be an interesting special case, much in analogy with the case of uniformly expanding polynomial mappings in one variable. We say that a set  $A \subset \mathbf{C}^2$  is a *hyperbolic set* for a mapping  $f : \mathbf{C}^2 \rightarrow \mathbf{C}^2$  if there are constants  $C < \infty$  and  $\lambda < 1$  and for each  $z \in A$  there is a splitting of the tangent bundle into stable and unstable directions, i.e.  $\mathbf{C}^2_z = E_z^s \oplus E_z^u$ , and

$$|Df^n(z)v| \leq C\lambda^n|v| \quad \text{for } z \in E_z^s$$

$$|Df^{-n}(z)v| \leq C\lambda^n|v| \quad \text{for } z \in E_z^u.$$

For our class of mappings satisfying (8), we will say that  $f$  is *hyperbolic* if  $J$  is a hyperbolic set for  $f$ . Not all of these mappings are hyperbolic, and it would be of interest to have criteria for hyperbolicity. Examples of hyperbolic maps are given by maps of the form (4) if  $p$  is uniformly expanding on the Julia set, and  $|a|$  is small. These mappings are rather easily seen to be hyperbolic (see [HO, FS], and [DN] for the real case). The only examples known at the moment are for  $|a|$  small or  $|c|$  large. See [HO] and [FS] for further properties of these maps.

The choice of  $J$  as the set on which  $f$  should have a hyperbolic splitting is justified by the following result [BS3, Theorem 6].

**Theorem.** *The following are equivalent:*

- (i)  $f$  has a hyperbolic splitting over the chain recurrent set.
- (ii)  $f$  has a hyperbolic splitting over the nonwandering set.
- (iii)  $f$  has a hyperbolic splitting over  $J$ .

Using elementary properties of hyperbolicity and the fact that the iterates of  $f$  are a normal family on  $\text{int } K^+$ , we obtain

**Theorem** [BS1]. *If  $f$  is hyperbolic, then the interior of  $K^+$  consists of the basins of finitely many hyperbolic sink orbits  $\{s_1, \dots, s_k\}$ .*

One of the reasons for considering the case of hyperbolic mappings is that stable manifolds always exist. The following holds for a smooth (not necessarily holomorphic) hyperbolic mapping.

**Stable Manifold Theorem.** *Let  $A$  be a compact hyperbolic set for a smooth diffeomorphism  $f$ . For every point  $x \in A$ ,  $W^s(x)$  is an immersed submanifold of dimension equal to that of  $E^s$ . Further,  $T_x W^s(x) = E_x^s$ , and the intersection of  $W^s(x)$  and  $W^u(x)$  at  $x$  is transverse.*

In the holomorphic case, we obtain more.

**Theorem.** *If  $f$  is hyperbolic, then the leaves of the stable and unstable foliations are complex submanifolds which are biholomorphically equivalent to  $\mathbf{C}$ .*

The relation between  $J^\pm$  and the stable/unstable foliations is as follows. It is not hard to show that  $W^s(J) = \bigcup_{x \in J} W^s(x) \subset J^+$ , and similarly for  $J^-$ . We even have:

**Theorem [BS1].** *If  $f$  is hyperbolic and  $|\det Df| \leq 1$ , then  $W^s(J) = J^+$ . If  $s_1, s_2, \dots, s_k$  are the sinks of  $f$  then  $W^u(J) = J^- - \{s_1, \dots, s_k\}$ .*

The unstable foliation cannot be extended through the sinks. In fact, Fornaess and Sibony [FS] show that there is no germ of a variety  $V_j$  with  $s_j \in V_j \subset J^-$ .

We let  $\mathcal{F}^s$  ( $\mathcal{F}^u$ ) denote the stable (unstable) foliation, and we describe what the local situation looks like. We may cover  $J^+$  with open sets  $U \subset J^+$  such that there is a coordinate system  $(x, y)$  such that  $U = J^+ \cap \{|x|, |y| < 1\}$ , and each leaf of  $\mathcal{F}^s \cap U$  is a graph  $\{y = \varphi(x) : |x| < 1\}$  for an analytic function  $\varphi$ . If we let  $T$  denote a transversal to the leaves of  $\mathcal{F}^s \cap U$ , then the leaves are parametrized by the set  $E := T \cap U$ . We note that with this choice of transversal, the measure  $\mu^+|_T$  gives a measure on the space of leaves of  $\mathcal{F}^s \cap U$ . If  $T'$  and  $T''$  are two transversals, then we may describe the space of leaves as  $E'$  and  $E''$ ; and we have a homeomorphism  $\chi : E' \rightarrow E''$ , which is given by following a leaf from  $T'$  to  $T''$ . Following Ruelle and Sullivan [RS], we say that the family of measures  $\{\mu^+|_T\}$  defines a *transversal measure* on  $\mathcal{F}^s \cap U$  if it is consistent with the homeomorphisms  $\{\chi\}$ , i.e. if  $\chi_* \mu^+|_{T'} = \mu^+|_{T''}$ . In [BS1, Theorem 6.5] it is shown that when  $f$  is hyperbolic, then the family  $\{\mu^+|_T\}$  does define a transversal measure on  $\mathcal{F}^s$ . This yields the following structure for  $\mu^+$  restricted to  $U$ , which may be written as  $\mu^+ \llcorner U$ :

$$\mu^+ \llcorner U = \int_{t \in E} [M_t] \mu^+|_T(dt),$$

where  $[M_t] = [\{y = \varphi_t(x)\}]$  denotes the current of integration over the leaf of  $\mathcal{F}^s \cap U$  passing through the point  $t \in E \subset T$ .

Using this local representation for  $\mu^\pm$  and the fact that the wedge product of currents of integration gives the current of integration over the intersection, i.e.  $[M^s] \wedge [M^u] = [M^s \cap M^u]$ , we obtain

$$\mu^+ \wedge \mu^- = \int_{t' \in E^s} \mu^+|_{T^s}(dt') \int_{t'' \in E^u} \mu^-|_{T^u}(dt'') [M_{t'}^s \cap M_{t''}^u],$$

which gives a sort of product structure to the measure  $\mu$ .

As a consequence, we obtain the result that *if  $f$  is hyperbolic, then the support of  $\mu$  is  $J$* . This property of the support of  $\mu$  yields the existence of periodic points:

**Corollary.** *If  $f$  is hyperbolic, then the periodic points for  $f$  are dense in  $J$ .*

Thus the conjecture of J.H. Hubbard that the periodic points for  $f$  are dense in  $J$  is verified in the special case of hyperbolic mappings.

Further applications of the theory of hyperbolic mappings yield the following (see [BS1, Corollary 7.9])  $\mu$  is mixing,  $\mu$  is the unique measure of maximal entropy, and  $\mu$  describes the distribution of periodic points. In fact mixing has recently been shown to hold for all maps (see [BS4]). It would be interesting know whether the other two properties continue to hold for more general mappings.

## References

- [Be] Bedford, E.: Survey of pluri-potential theory. Several complex variables. Proceedings of the Mittag-Leffler-Institute, 1987–88, J.-E. Fornaess (ed.), Princeton Univ. Press, 1991
- [BS1] Bedford, E., Smillie, J.: Polynomial diffeomorphisms of  $C^2$ : Currents, equilibrium measure and hyperbolicity. Invent. math. **103** (1991) 69–99
- [BS2] Bedford, E., Smillie, J.: Fatou-Bieberbach domains arising from polynomial automorphisms. Indiana Univ. Math. J. (to appear)
- [BS3] Bedford, E., Smillie, J.: Polynomial diffeomorphisms of  $C^2$ : Stable manifolds and recurrence. Preprint
- [BS4] Bedford, E., Smillie, J.: Polynomial diffeomorphisms of  $C^2$ : Ergodicity, exponents and entropy of the equilibrium measure. Preprint
- [BT] Bedford, E., Taylor, B.A.: The Dirichlet problem for a complex Monge-Ampère equation. Invent. math. **50** (1976) 543–571
- [BC] Benedicks, M., Carleson, L.: The dynamics of the Hénon map. Preprint
- [Bi] Bielefeld, B.: Conformal dynamics problem list
- [Br] Brolin, H.: Invariant sets under iteration of rational functions. Ark. Mat. **6** (1965) 103–144
- [C1] Cegrell, U.: On the discontinuity of the complex Monge-Ampère operator. Proceedings, Toulouse 1983. (Lecture Notes in Mathematics, vol. 1094.) Springer, Berlin Heidelberg New York
- [C2] Cegrell, U.: Discontinuité de l'opérateur de Monge-Ampère complexe. C.R. Acad. Sci. Paris. Sér. I Math. **296** (1983) 869–871
- [DN] Devaney, R., Nitecki, Z.: Shift automorphisms in the Hénon mapping. Comm. Math. Phys. **67** (1979) 137–146
- [DE] Dixon, P.G., Esterle, J.: Michael's problem and the Poincaré-Fatou-Bieberbach phenomenon. Bull. AMS **15** (1986) 127–187
- [FM] Friedland, S., Milnor, J.: Dynamical properties of plane polynomial automorphisms. Ergodic Theory Dyn. Syst. **9** (1989) 67–99
- [FS] Fornæss, J.-E., Sibony, N.: Complex Hénon mappings in  $C^2$  and Fatou Bieberbach domains. Preprint
- [Hn1] Hénon, M.: A two-dimensional mapping with a strange attractor. Comm. Math. Phys. **50** (1976) 69–77
- [Hn2] Hénon, M.: Numerical study of quadratic area preserving mappings, Q. Appl. Math. **27** (1969) 291–312
- [Hr] Herman, M.: Recent results and some open questions on Siegel's linearization theorem of germs of complex analytic diffeomorphisms of  $C^n$  near a fixed point. VIIth International Congress on Mathematical Physics (Marseille, 1986). World Scientific Publishing, Singapore 1987, pp. 138–184
- [Hu] Hubbard, J.: The Hénon mapping in the complex domain. In: M. Barnsley and S. Demko (eds.), Chaotic Dynamics and Fractals. Academic Press, New York 1986, pp. 101–111
- [HO] Hubbard, J., Oberste-Vorth, W.: Hénon mappings in the complex domain. In preparation

- [K] Katok, A.: Nonuniform hyperbolicity and structure of smooth dynamical systems. Proceedings of the International Congress of Mathematicians, Warsaw 1983
- [L1] Lelong, P.: Fonctions plurisousharmoniques et formes différentielles positives. Gordon and Breach, New York 1968
- [L2] Lelong, P.: Discontinuité et annulation de l'opérateur de Monge-Ampère complexe. Séminaire P. Lelong-P. Dolbeault-H. Skoda (1983). (Lecture Notes in Mathematics, vol. 1028.) Springer, Berlin Heidelberg New York 1983, pp. 219–224
- [Lo] Lopes, A.: Equilibrium measures for rational maps. *Ergod. Theor. Dyn. Syst.* **6** (1986) 393–399
- [M] Milnor, J.: Non-expansive Hénon maps. *Adv. Math.* **69** (1988) 109–114
- [O] Oberste-Vorth, R.: Complex horseshoes and the dynamics of mappings of two complex variables. Thesis, Cornell University 1987
- [RR] Rosay, J.-P., Rudin, W.: Holomorphic maps from  $\mathbf{C}^n$  to  $\mathbf{C}^n$ . *Trans. AMS* **310** (1988) 47–86
- [RS] Ruelle, D., Sullivan, D.: Currents, flows, and diffeomorphisms. *Topology* **14** (1975) 319–327
- [S] Shub, M.: Global Stability of Mappings. Springer, Berlin Heidelberg New York 1987
- [Sm] Smillie, J.: The entropy of polynomial diffeomorphisms of  $\mathbf{C}^2$ . *Ergod. Theor. Dyn. Syst.* **10** (1990) 823–827
- [Su1] Sullivan, D.: Cycles for the dynamical study of foliated manifolds and complex manifolds. *Invent. math.* **36** (1976) 225–255
- [Su2] Sullivan, D.: Quasiconformal homeomorphisms and dynamics II
- [T] Tortrat, P.: Aspects potentialistes de l'itération des polynômes. Séminaire de Théorie du Potentiel, Paris, no. 8. (Lecture Notes in Mathematics, vol. 1235.) Springer, Berlin Heidelberg New York 1987
- [U1] Ueda, T.: Local structure of analytic transformations of two complex variables. *J. Math. Kyoto Univ.* **26** (1986) 233–261
- [U2] Ueda, T.: Local structure of analytic transformations of two complex variables. II. Preprint
- [Z] Zehnder, E.: A simple proof of a generalization of a Theorem by C.L. Siegel, J. Palis & M. do Carmo (eds.) *Geometry and Topology*. (Lecture Notes of Mathematics, vol. 597.) Springer, Berlin Heidelberg New York 1977, pp. 855–866

# Precise Analysis of $\bar{\partial}_b$ and $\bar{\partial}$ on Domains of Finite Type in $\mathbb{C}^2$

Michael Christ \*

Department of Mathematics, University of California, 405 Hilgard Ave.  
Los Angeles, CA 90024, USA

## 0. Introduction

In the past four years or so marked progress has been achieved in understanding precise regularity properties of solutions of the  $\bar{\partial}$  equation on smoothly bounded two-dimensional complex manifolds, and of  $\bar{\partial}_b$  on three-dimensional CR manifolds, under hypotheses of pseudoconvexity and finite type. Hölder and  $L^p$  Sobolev estimates, sharp nonisotropic estimates in  $L^2$  norms, and pointwise bounds for Szegö and Bergman kernels have been obtained in partially overlapping works by a number of authors. This article represents a summary of those developments, and of two applications.

The approach taken is that of canonical solutions and  $L^2$  theory. Important developments have also occurred in the method of explicit solution via integral kernels, but are neglected here because of the author's lack of expertise. See [Be, Fo, Ra].

Some notation:  $L_s^p$  denotes the Sobolev space of all functions having  $s$  derivatives in  $L^p$ , for  $1 < p < \infty$  and  $s \geq 0$ .  $A_\alpha$  denotes the space of functions Hölder continuous of order  $\alpha > 0$ , with the usual convention when  $\alpha$  is an integer [St].

## 1. Definitions

The setting for our first group of results is a  $C^\infty$ , compact (real) three-dimensional manifold  $M$  without boundary. A CR structure on such a manifold is a smooth, complex one-dimensional subbundle  $T^{0,1}$  of the complexified tangent bundle of  $M$ , satisfying  $T_x^{0,1} \cap T_x^{1,0} = \{0\}$  for all  $x \in M$ , where  $T^{1,0} = \overline{T^{0,1}}$ . Let  $B^{0,1}$  be the bundle dual to  $T^{0,1}$ . Canonically associated to the CR structure is a first-order partial differential operator  $\bar{\partial}_b$ , mapping functions to sections of  $B^{0,1}$ , defined by

$$\bar{\partial}_b(f)(v) = v(f)$$

for  $v \in T_x^{0,1}$  for some  $x$ . Thus  $\bar{\partial}_b f$  is a portion of the usual differential  $df$ .

\* Alfred P. Sloan fellow. Research also supported by the Institut des Hautes Etudes Scientifiques and National Science Foundation.

Although the problems to be considered here have also a global aspect, the bulk of the analysis will take place in a small coordinate patch, about an arbitrary point  $x_0$ . Then  $B^{0,1}$  trivializes and  $\bar{\partial}_b$  may be regarded as a complex vector field, which we write as  $\bar{\partial}_b = X + iY$ , where  $X, Y$  are smooth real vector fields, linearly independent at every point. Since  $M$  has dimension three,  $\bar{\partial}_b$  is not elliptic.

Fix a real vector field  $T$  which is linearly independent of  $X, Y$  at  $X_0$ . Define the Levi form  $\lambda$  as a real-valued function by

$$[X, Y](x) = \lambda(x)T + O(X, Y).$$

**Definition 1.1.**  *$M$  is said to be pseudoconvex at  $x_0$  if  $\lambda$  does not change sign in some neighborhood of  $x_0$ .*

The usual distinction between pseudoconvexity and pseudoconcavity is lost here; by replacing  $T$  by  $-T$  if necessary, we may assume henceforth that  $\lambda \geq 0$ .

**Definition 1.2.**  *$M$  is said to be of finite type at  $x_0$  if  $X, Y$  satisfy the Hörmander condition that they, together with all their Lie brackets of all orders, should span the tangent space to  $M$  at  $x_0$ .*

$M$  is said to be pseudoconvex (respectively of finite type) if it is pseudoconvex (respectively of finite type) at every point. Both definitions are independent of various arbitrary choices involved – for instance the choice of  $T$ .  $x_0$  is a point of type  $m$  if commutators of length  $m$ , but no smaller length, suffice to span the tangent space. ( $[X, Y]$  is said to have length two,  $[X, [X, Y]]$  length three, and so on.)  $M$  of type  $m = \max_{x \in M} \text{type}(x)$ . An equivalent condition is that some differential monomial  $D$  in  $X, Y$  should satisfy  $D\lambda(x_0) \neq 0$ .

The principal examples of CR manifolds are boundaries of open domains in  $\mathbb{C}^2$ ;  $T_x^{0,1}$  is taken to be the space of anti-holomorphic tangent vectors to  $\mathbb{C}^2$ , which are tangent to  $M$  at  $x$ . Globally real analytic boundaries are always of finite type. In this case  $\bar{\partial}_b f$  may be computed by extending  $f$  smoothly to  $\mathbb{C}^2$ , applying  $\bar{\partial}$  to the extension, and restricting the resulting  $(0, 1)$  form to  $M$ . Thus  $\bar{\partial}_b$  annihilates the restriction to  $M$  of any holomorphic function, hence has globally an infinite-dimensional kernel. Consequently  $\bar{\partial}_b$  cannot satisfy any subelliptic estimate, despite the Hörmander condition.

It is necessary to choose a particular solution of  $\bar{\partial}_b u = f$  in order to have any regularity theory; one candidate is the solution of minimal norm (after fixing a Hermitian metric on  $B^{0,1}$  and a volume form on  $M$ ) [K3]. Thus one studies

$$\bar{\partial}_b u = f \quad \text{with} \quad u \perp \text{Kernel}(\bar{\partial}_b).$$

## 2. $L^2$ Theory

Having fixed a metric and volume form, we may speak of  $L^2$  sections of  $B^{0,1}$ ; henceforth we shall use the symbol  $L^2$  to refer either to sections or to functions. In order to have any satisfactory existence theory, one requires

**Hypothesis 2.1.**  $\bar{\partial}_b$  is assumed to have closed range in  $L^2$ .

More explicitly, if  $f \in L^2$  is in the closure of  $\bar{\partial}_b(L^2) \cap L^2$ , then  $f = \bar{\partial}_b u$  for some  $u \in L^2$ , and upon choosing the solution  $u$  orthogonal to the kernel, one has  $\|u\| \leq C\|f\|^1$ . This hypothesis is violated generically; see §7. However:

**Theorem 2.2** [K2]. *Suppose that  $M$  is pseudoconvex and may be realized as the boundary (in the  $C^\infty$  sense) of a complex variety in some  $\mathbb{C}^n$ . Then  $\bar{\partial}_b$  has closed range in  $L^2(M)$ .*

For the case of the boundary of an open domain see [BS].

The principal result in the  $L^2$  theory is then

**Theorem 2.3** [K3]. *Suppose that  $M$  is a three-dimensional CR manifold, pseudoconvex and of finite type  $\leq m$ , on which  $\bar{\partial}_b$  has closed range in  $L^2$ . Then for any  $C^\infty$  functions  $\varphi, \psi$  satisfying  $\varphi\psi \equiv \psi$ , for any  $s \geq 0$  and any  $f \in L^2$ , the unique solution  $u$  of  $\bar{\partial}_b u = f$  orthogonal to Kernel( $\bar{\partial}_b$ ) satisfies*

$$\|\psi u\|_{L^2_{s+\delta}} \lesssim \|\varphi f\|_{L^2_s} + \|f\|.$$

In this theorem and the next,  $\delta = m^{-1}$ .

Indeed being orthogonal to the kernel,  $u = \bar{\partial}_b^* v$  where  $\|v\| \lesssim \|f\|$ , so this follows from the fundamental local result:

**Theorem 2.4.** *Suppose that  $U \Subset U'$  are open sets and that  $\bar{\partial}_b \bar{\partial}_b^* v = f$  on  $U'$ . Then*

$$\|\bar{\partial}_b^* v\|_{L^2_{s+\delta}(U)} \lesssim \|f\|_{L^2_s(U')} + \|v\|_{L^2(U')}.$$

A well-known analogue is that  $(X^2 + Y^2)v = f$  implies the same conclusion for  $Xv$  and  $Yv$ .

The fundamental innovation in the proofs of Theorems 2.2 and 2.3 was a rather simple microlocal analysis<sup>2</sup>. In brief: Choose coordinates  $(x, y, t)$  in which  $x_0 = 0$ ,  $X(0) = \partial/\partial x$ ,  $Y(0) = \partial/\partial y$ , and  $T \equiv \partial/\partial t$ . Let  $(\xi', \tau) \in \mathbb{R}^2 \times \mathbb{R}$  be dual variables. Near 0,  $\bar{\partial}_b$  is elliptic where  $|\tau| \leq C|\xi'|$ . Let  $P^+, P^-$  be classical pseudodifferential operators of order zero, with symbols supported where  $\tau > 0$  and  $|\tau| \geq C|\xi'|$ , or where  $\tau < 0$  and  $|\tau| \geq C|\xi'|$ , respectively.

<sup>1</sup>  $\|\cdot\|$ , with no subscript, will always denote the  $L^2$  norm.

<sup>2</sup> The same idea is found in [HN], in the context of  $\square_b$ .

$$\|\bar{\partial}_b P^- u\|^2 \gtrsim \|XP^- u\|^2 + \|YP^- u\|^2 - C\|u\|^2,$$

modulo some error terms, by Gårding's inequality. The right-hand side controls  $\|u\|_{L^2_\delta}^2 - C\|u\|^2$ , by the result of [H] and [K1], so  $\bar{\partial}_b$  is subelliptic in the support of the symbol of  $P^-$ .  $\bar{\partial}_b$  is not subelliptic in the support of the symbol of  $P^+$ , but a similar application of Gårding's inequality demonstrates that instead,  $\bar{\partial}_b \bar{\partial}_b^*$  is subelliptic there. Thus one studies  $P^+ v$ , where  $u = \bar{\partial}_b^* v$ . One cannot eliminate the microlocalization by simply analyzing the equation  $\bar{\partial}_b \bar{\partial}_b^* v = f$ , as  $\bar{\partial}_b \bar{\partial}_b^*$  fails to be subelliptic in the support of  $P^-$ .

### 3. Underlying Geometry

Of fundamental importance in the refined analysis of  $\bar{\partial}_b$  is a second geometric structure induced on  $M$  by the CR structure [FS, RS, NSW, Sa, FP]. A curve  $\gamma : [0, r] \mapsto M$  is said to be admissible if  $\gamma$  is absolutely continuous,  $|\gamma'(t)| \leq 1$  for almost every  $t$ , and if for almost every  $t$ ,  $\gamma'(t) \in T^{0,1} \oplus T^{1,0}$ . Henceforth  $M$  is always assumed to be connected. Define:

$$\varrho(x, y) = \inf\{r : \exists \text{ an admissible curve } \gamma : [0, r] \mapsto M \text{ with } \gamma(0) = x \text{ and } \gamma(r) = y\}. \quad (3.1)$$

It is an easy theorem that under the hypothesis of finite type, any two points can be joined by an admissible curve. Then  $\varrho$  becomes a metric, in the sense of point-set topology. However it is not a Riemannian metric, and  $\varrho(x, y)$  may be as large as  $c[\text{dist}(x, y)]^{1/m}$ , where  $\text{dist}$  is a Riemannian distance and  $M$  is of type  $m$ .

Let  $B(x, r)$  denote the open ball centered at  $x$ , with respect to  $\varrho$ . Let  $B_0$  denote the unit ball in  $\mathbb{R}^3$ . Given  $x \in M$  and  $r > 0$ , define  $\Phi : B_0 \mapsto M$  by

$$\Phi_{x,r}(u) = \exp(ru_1 X + ru_2 Y + \delta(x, r)u_3 T)(x)$$

where

$$\delta(x, r) = \sum_{0 \leq |D| \leq m-2} r^{|D|+2} |D\lambda(x)| \quad (3.2)$$

and  $D$  ranges over all differential monomials in  $X, Y$ , with  $|D|$  defined to be the number of factors of  $X, Y$ . Then  $\Phi_{x,r}$  is a diffeomorphism for all small  $r$ . Define  $\tilde{B}(x, r) = \Phi_{x,r}(B_0) \subset M$ . Then there exists  $c$  such that  $B(x, cr) \subset \tilde{B}(x, r) \subset B(x, c^{-1}r)$  for all  $x, r$ .

**Proposition 3.1** [NSW].  $|\tilde{B}(x, r)| \sim r^2 \delta(x, r)$ .

**Corollary 3.2.** *The ordered triple  $(M, \varrho, \text{volume form})$  is a space of homogeneous type, in the sense of the theory of singular integral operators.*

See [CW] for the definition, and also [Ch7] for a recent exposition.

Fixing  $x, r$ , pull  $X, Y$  back to  $B_0$  by

$$(\hat{X}f)(u) = r \cdot (\check{X}f)(\Phi_{x,r}u)$$

where  $\check{f} = f \circ \Phi_{x,r}^{-1}$ , and similarly for  $\hat{Y}$ . Pull  $\bar{\partial}_b$  back in the same way.

**Theorem 3.3 [NSW].**

- The coefficients of  $\hat{X}, \hat{Y}$  are  $C^\infty$  functions of  $u \in B_0$ , uniformly in  $x, r$ .
- The CR structure on  $B_0$  associated to  $(\bar{\partial}_b)^\wedge$  is of finite type, uniformly in  $x, r$ .

The last assertion means that  $\hat{X}, \hat{Y}$  satisfy the Hörmander condition in a uniform way, hence satisfy a subelliptic estimate completely independent of  $x, r$ . Thus after pulling back to  $B_0$ , one has both upper and lower bounds on  $(\bar{\partial}_b)^\wedge$ , uniformly in  $x, r$ ; the “rescaling” maps  $\Phi_{x,r}$  thus provide a way to analyze  $\bar{\partial}_b$  uniformly at all points and all scales, and the geometry defined by  $\varrho$  is precisely adapted to the study of operators such as  $\bar{\partial}_b$  and  $X^2 + Y^2$ .

Define  $\tilde{\nabla}f = (Xf, Yf)$  and define  $\tilde{\nabla}^\alpha$  analogously. Combining Theorems 2.4 and 3.3 yields

**Corollary 3.4.** *For any  $N$  there exists  $N'$  with the following property: for any  $x \in M$ , any  $0 < r \leq 1$ , and any  $v, g$  satisfying*

$$\begin{aligned} \|\tilde{\nabla}^\alpha g\|_{L^2(\tilde{B}(x,r))} &\leq r^{-|\alpha|} \quad \forall |\alpha| \leq N', \\ \bar{\partial}_b \bar{\partial}_b^* v &= g \quad \text{in } \tilde{B}(x,r), \\ \|\bar{\partial}_b^* v\|_{L^2(\tilde{B}(x,r))} &\leq r, \text{ and} \\ \|v\|_{L^2(\tilde{B}(x,r))} &\leq r^2, \end{aligned}$$

one has  $\bar{\partial}_b^* v \in C^\infty(\tilde{B}(x,r))$  and

$$\|\tilde{\nabla}^\alpha \bar{\partial}_b^* v\|_{L^\infty(\tilde{B}(x,r/2))} \leq C_N r^{1-|\alpha|} |B(x,r)|^{-1/2} \quad \forall |\alpha| \leq N.$$

Let it be stressed that  $C_N$  is independent of  $x, r$ . This result is one of the principal ingredients in the proof of Theorem 4.1 below, and explains why that theorem should hold.

## 4. Main Results for CR Manifolds

Throughout this section  $M$  is assumed to be a smooth, compact three-dimensional CR manifold without boundary, pseudoconvex and of finite type not exceeding  $m$ , and  $\bar{\partial}_b$  is assumed to have closed range. In addition to regularity results for  $\bar{\partial}_b$ , a goal has been to elucidate the nature of certain operators associated to  $\bar{\partial}_b$ , as had been done previously in the strictly pseudoconvex case [FS, Fe1, BSj]. These are principally the Szegö projection, that is, the orthogonal projection of  $L^2(M)$  onto the kernel of  $\bar{\partial}_b$ , and the relative solving operator  $G$ , which is bounded and linear on  $L^2$  and is uniquely defined by the relations  $\bar{\partial}_b Gf = \pi f$  and  $Gf \perp \text{Kernel}(\bar{\partial}_b)$ , where  $\pi f$  denotes the orthogonal projection of  $f$  onto the range of  $\bar{\partial}_b$ .<sup>3</sup> Denote

---

<sup>3</sup> It can be shown that the kernel and cokernel of  $\bar{\partial}_b$  both have infinite dimension.

by  $S(x, y), G(x, y)$  their distribution-kernels, which are  $C^\infty$  off of the diagonal as a consequence of Theorem 2.3.

Define  $V(x, y) = |B(x, \varrho(x, y))|$ . Let  $D$  denote any differential monomial in  $X, Y$ , acting in either or both of the variables  $x, y$ , and let  $\ell$  be the number of factors of  $X, Y$ .

**Theorem 4.1** [Ch2, FK2, Ch5].  $|DG(x, y)| \lesssim \varrho(x, y)^{1-\ell} V(x, y)^{-1} \quad \forall x \neq y.$

On the level of operators  $S = I - G \circ \bar{\partial}_b$ , so one has also:

**Corollary 4.2** [Ch1, NRSW2, FK2].  $|DS(x, y)| \lesssim \varrho(x, y)^{-\ell} V(x, y)^{-1} \quad \forall x \neq y.$

The estimates for  $S$  and its derivatives, together with its  $L^2$  boundedness, imply that it is a singular integral operator in the sense of the Calderón-Zygmund theory, with respect to the geometric structure on  $M$  imposed by  $\varrho$  and the volume form. One of the fundamental results of that theory is that  $L^2$  boundedness implies  $L^p$  boundedness.

In the same way,  $L^p$  boundedness of  $X \circ G, Y \circ G$  would follow from  $L^2$  boundedness. A second fundamental result of singular integral theory is the  $T(1)$  theorem [DJ], which provides a necessary and sufficient condition for  $L^2$  boundedness, and whose hypotheses are often checkable in practice. In the present instance the main condition to be verified is “weak boundedness”, and this follows from some of the ingredients of the proof of the pointwise estimates for the kernels. Thus one has the case  $s = 0$  of:

**Theorem 4.3** [Ch1]. *Let  $p \in (1, \infty)$  and  $s \geq 0$ . Then  $X \circ G$  and  $Y \circ G$  are bounded from  $L_s^p$  to  $L_s^p$ . Moreover  $G : L_s^p \mapsto L_{s+m-1}^p$ , and  $S : L_s^p \mapsto L_s^p$ .*

For  $s > 0$  one needs also to commute differential operators past  $G$  and  $S$ ; this may be accomplished via microlocal arguments using standard pseudodifferential operators as in [FK1, Ch1, Ch5] or by the formalism of [CNS2]. Singular and fractional integral operators such as  $S, G$  are studied systematically in [NRSW2, CNS2].

**Corollary 4.4** [Ch1, Ch5]. *If  $U \Subset U' \subset M$  are open,  $\bar{\partial}_b u = f$  and  $u = \bar{\partial}_b^* v$  in  $U'$  for some  $v \in L^2$ , and if  $f \in L_s^p(U')$ , then  $u \in L_{s+m-1}^p(U)$ .*

There are analogous results in the scale of Hölder spaces:

**Theorem 4.5** [FK1]. *If  $U \Subset U' \subset M$  are open,  $\bar{\partial}_b u = f$  and  $u = \bar{\partial}_b^* v$  in  $U'$  for some  $v \in L^2$ , if  $\alpha > 0$  and if  $f \in A_\alpha(U')$ , then  $u \in A_{\alpha+m-1}(U)$ <sup>4</sup>.*

<sup>4</sup> This is proved with a nonstandard definition of  $A_\alpha$  for  $\alpha \in \mathbb{Z}$  in [FK1], but follows with the usual definition [St] from the machinery of [NRSW2, CNS2].

Another fundamental operator is  $\square_b^1 = \bar{\partial}_b \bar{\partial}_b^*$ ; associated to it is a relative solving operator whose distribution-kernel satisfies estimates like those for  $G$ , but with a power of  $r^{2-\ell}$  [CNS2]. It has mapping properties like those of  $G$ , but gains two derivatives in the  $X, Y$  directions instead of one.

Basic difficulties in the proofs of these results are:

- There is no standard pseudodifferential calculus which permits inversion (modulo the kernel) of operators such as  $\bar{\partial}_b$ ; one is led into classes  $S_{\varrho,\delta}^k$  with  $\varrho < \delta$ .
- There is no single model by which one may approximate, as by the Heisenberg group in the strictly pseudoconvex case; and the natural model hypersurfaces  $\{\text{Im}(z_2) = P(z_1)\}$  where  $P$  is a subharmonic, nonharmonic polynomial seem practically as difficult to analyze as the general case.
- The nonexistence of good holomorphic support functions [KN] renders the construction of explicit parametrices difficult.
- The lifting technique of [RS, NSW] and [Sa] is not (directly) applicable.
- The operators  $P^\pm$  of the microlocal analysis sketched earlier do not preserve the type of estimates sought, thus introducing spurious singularities.<sup>5</sup>

The first proofs of these results were rather involved, but subsequently substantial simplification has occurred. See [Ch5].

## 5. The $\bar{\partial}$ -Neumann Problem

Let  $\Omega \subset \mathbb{C}^2$  be pseudoconvex and smoothly bounded.  $\Omega$  is said to be of finite type if  $\partial\Omega$  is of finite type in the sense already defined; an equivalent notion of finite type (in  $\mathbb{C}^2$ ) is that  $\partial\Omega$  should have a bounded order of contact with all complex curves.

Fix a Riemannian metric on a neighborhood of  $\bar{\Omega}$ , let  $r$  be the signed geodesic distance to  $\partial\Omega$ , and set  $\square = \bar{\partial}\bar{\partial}^* + \bar{\partial}^*\bar{\partial}$ , on  $(0, 1)$  forms. The  $\bar{\partial}$ -Neumann problem is the boundary value problem

$$\square u = f \text{ on } \Omega \text{ with } \begin{cases} u \lrcorner \bar{\partial}r = 0 & \text{on } \partial\Omega \\ \bar{\partial}u \lrcorner \bar{\partial}r = 0 & \text{on } \partial\Omega \end{cases}$$

for  $(0, 1)$  forms on  $\Omega$ , where  $\lrcorner$  denotes the interior product of forms with respect to the Riemannian metric.  $\square$  is an elliptic second-order system, but the boundary conditions are non-coercive. Let  $N$  be the Neumann operator, which inverts  $\square$  on  $(0, 1)$  forms, modulo its finite-dimensional kernel.  $N$  sends  $L_s^2(\Omega)$  to  $L_{s+(2/m)}^2(\Omega)$  [K1, RS] for all  $s \geq 0$ , assuming  $\partial\Omega$  to be of finite type  $m$ .

Precise results on the regularity properties of  $N$  in various function spaces have been obtained by [CNS2], generalizing previously known results for the strictly pseudoconvex case [GS]. One of the main steps is a computation of  $N$  in terms of simpler operators, as follows. Let  $\bar{\omega}_1, \bar{\omega}_2$  be  $C^\infty$  and constitute an orthonormal basis for the space of  $(0, 1)$  forms at each point in a neighborhood

<sup>5</sup> This is particularly evident in [Ch1] and [FK1].

of  $\partial\Omega$ , with  $\bar{\omega}_2 = \sqrt{2}\bar{\partial}r$ . Let  $\{\bar{L}_1, \bar{L}_2\}$  be a dual basis of anti-holomorphic vector fields. Let  $G$  solve for  $(0, 1)$  forms  $f$ :

$$\square Gf = f \text{ on } \Omega \quad \text{with } Gf = 0 \text{ on } \partial\Omega ;$$

write  $Gf = G_1 f + G_2 f$  where  $G_i f$  is a scalar multiple of  $\bar{\omega}_i$ . Let  $R$  denote restriction to  $\partial\Omega$ . Let  $P$  solve for  $(0, 1)$  forms  $f$  defined on  $\partial\Omega$  (that is,  $f$  is a scalar multiple of  $\bar{\omega}_1$ )

$$\begin{cases} \square Pf = 0 & \text{on } \Omega \text{ to infinite order at } \partial\Omega \\ Pf \equiv f & \text{on } \partial\Omega. \end{cases}$$

Finally let  $K$  denote the relative solution operator for  $\bar{\partial}_b \bar{\partial}_b^* = \square_b^1$  on  $\partial\Omega$ . By a  $\psi do$  of order  $n$  we shall mean a classical pseudodifferential operator with symbol in  $S_{1,0}^n$ .

**Theorem 5.1** [CNS2].  $N = G + P \circ (\square^- K \Gamma^+ + Q) \circ R \circ \bar{L}_2 G_1$  plus lower order terms.

Here

- $Q$  is a  $\psi do$  on  $\partial\Omega$  of order  $-1$
- $\Gamma^+$  is a  $\psi do$  of order  $0$  which microlocalizes to the region in phase space in which  $\bar{\partial}_b$  is not subelliptic
- $\square^-$  is a  $\psi do$  of order  $+1$  which may be explicitly computed modulo a symbol of order  $-1$ .

The “lower order” terms are of a form similar to the principal one, so that in practice, any mapping property which one can establish for the main term, may also be proved for the lower-order part. The only ingredient which is not computable and quite well understood is  $K$ , and one must make do with properties established for it using the results of §4.

Combining Theorem 5.1 with results on  $K$ , plus commutation properties of  $K$  with vector fields, results in

**Theorem 5.2** [CNS2]. Let  $q$  be a polynomial of degree at most  $2$  in  $L_1, \bar{L}_1$ . Then

$$\begin{aligned} q(L_1, \bar{L}_1) \circ N : L_s^p(\Omega) &\mapsto L_s^p(\Omega) & \forall p \in (1, \infty), s \geq 0 \\ N : A_\alpha &\mapsto A_{\alpha + \frac{2}{m}} & \forall \alpha > 0. \end{aligned}$$

Also  $N : L^\infty \mapsto A_\alpha$  for all  $\alpha < \frac{2}{m}$ .

Let  $B : L^2(\Omega) \mapsto L^2(\Omega)$  be the Bergman projection onto the holomorphic functions. If  $\bar{\partial}f = 0$ , then  $\bar{\partial}(\bar{\partial}^* Nf) = f$ , so that  $\bar{\partial}^* N$  solves the  $\bar{\partial}$  equation and thus is of much interest.

**Corollary 5.3** [CNS2].

$$\begin{aligned}\bar{\partial}^* N, L_1 \circ \bar{\partial}^* N, \bar{L}_1 \circ \bar{\partial}^* N : L_s^p(\Omega) &\mapsto L_s^p(\Omega) \quad \forall p \in (1, \infty), s \geq 0 \\ \bar{\partial}^* N : A_\alpha(\Omega) &\mapsto A_{\alpha+m-1}(\Omega) \quad \forall \alpha > 0 \\ \bar{\partial}^* N : L^\infty(\Omega) &\mapsto A_\beta(\Omega) \quad \forall \beta < m^{-1} \\ B : L_s^p(\Omega) &\mapsto L_s^p(\Omega) \quad \forall p \in (1, \infty), s \geq 0.\end{aligned}$$

In order to formulate pointwise estimates on the Bergman kernel  $B(z, w)$ , extend  $\varrho$  to  $\bar{\Omega} \times \bar{\Omega}$  by:

$$\varrho(z, w) = \max \left( \varrho(\pi(z), \pi(w)), \mu(z), \mu(w) \right),$$

where  $\mu$  is defined by the relation

$$\delta(\pi(x), \mu(x)) = r(x), \tag{5.1}$$

and where  $\pi(x)$  denotes the point of  $\partial\Omega$  closest to  $x$ , and  $\delta$  is the invariant defined in (3.2). In the strictly pseudoconvex case,  $\mu \sim r^{1/2}$ , but on domains of finite type  $m$  one has only  $c_1 r^{1/2} \leq \mu \leq c_2 r^{1/m}$ . Let  $D$  be any differential monomial in  $L_1, \bar{L}_1, L_2, \bar{L}_2$ , with each factor permitted to act in either of the variables  $z, w$ , with  $k$  factors of  $L_1, \bar{L}_1$ , and  $\ell$  factors of  $L_2, \bar{L}_2$ .

**Theorem 5.4** [Mc, NRSW2].  $|DB(z, w)| \leq C_{k,\ell} \varrho(z, w)^{-2-k} \delta(\pi(z), \varrho(z, w))^{-2-\ell}$ .

## 6. Higher Dimensions

The question of Hölder and  $L^p$  estimates and the nature of the Szegö and Bergman kernels is at present poorly understood in higher dimensions; what is clear is that everything is more subtle. The notion of finite type is itself more recondite. A boundary point  $x$  is said [D] to be of finite type  $m$  if no one-dimensional complex variety in the ambient space has order of contact greater than  $m$  with the boundary at  $x$ .<sup>6</sup> It is a theorem that the set of boundary points of finite type is open [D]. Given pseudoconvexity, finite type is a necessary [C1] and sufficient [C2] condition for the  $\bar{\partial}$ -Neumann problem to satisfy a subelliptic estimate. The proofs of these assertions are several orders of magnitude more difficult than in the two-dimensional case.

**Question.** Let  $\Omega$  be pseudoconvex and of finite type. Does there exist  $\varepsilon > 0$  such that the Neumann operator maps  $A_\alpha(\bar{\Omega})$  to  $A_{\alpha+\varepsilon}(\bar{\Omega})$  and  $L^\infty$  to  $A_\varepsilon$ ?

The principal difficulty in higher dimensions stems from the nature of  $\square_b$  as a system, rather than a scalar equation. There has been definite progress in the simpler case where the Levi form of  $\Omega$  is diagonalizable [M1, M2, FKM].

<sup>6</sup> In  $\mathbb{C}^n$  for  $n > 2$ , this is not equivalent to the natural condition in terms of vector fields.

However certain negative as well as positive results have been obtained, so that even the diagonalizable case is more subtle than might have been anticipated.

In  $\mathbb{C}^3$  consider a domain  $\{\operatorname{Im}(z_3) > |z_1|^{2k} + |z_2|^2\}$ , for  $k = 2, 3, \dots$ . Identify the boundary with  $\mathbb{C}^2 \times \mathbb{R}$  via the map  $(z_1, z_2, t) \mapsto (z_1, z_2, t + i(|z_1|^{2k} + |z_2|^2))$ . Set

$$L_1 = \frac{\partial}{\partial z_1} + ik\bar{z}_1|z_1|^{2k-2} \frac{\partial}{\partial t}, \quad L_2 = \frac{\partial}{\partial z_2} + i\bar{z}_2 \frac{\partial}{\partial t}, \quad T = \frac{\partial}{\partial t}.$$

Then after being pulled back to  $\mathbb{C}^2 \times \mathbb{R}$ ,  $\square_b := \bar{\partial}_b \bar{\partial}_b^* + \bar{\partial}_b^* \bar{\partial}_b$  on  $(0, 1)$  forms becomes

$$\square_b = \begin{pmatrix} \mathcal{L} & 0 \\ 0 & \mathcal{L} \end{pmatrix} \quad \text{where } \mathcal{L} = -(\bar{L}_1 L_1 + L_2 \bar{L}_2).$$

Thus one need only analyze  $\mathcal{L}^{-1}$ .

What geometric structure should underlie the analysis? Write  $\bar{L}_j = X_j + iY_j$ . Let  $\varrho$  denote the metric defined in terms of  $X_1, Y_1, X_2, Y_2$  as in (3.1); one might hope that  $\mathcal{L}^{-1}$  should behave as an operator smoothing of order two with respect to the balls associated to  $\varrho$ , as for the relative fundamental solution for  $\square_b$  in the two-dimensional case. This is false. Let  $\lambda(z_1, z_2, t) = |z_1|^{2k-2}$  be the degenerate eigenvalue of the Levi form, up to a constant factor. Let  $\tilde{\varrho}$  be the metric associated to the weaker, but still subelliptic, second-order operator  $\tilde{\mathcal{L}} = -(X_1^2 + Y_1^2 + X_2\lambda X_2 + Y_2\lambda Y_2)$  according to the theory of [FP]. One has  $\tilde{\varrho} \geq \varrho$ , but no majorization  $\tilde{\varrho} \leq C\varrho$  near where  $\lambda$  vanishes. Denote by  $B_\varrho$  and  $B_{\tilde{\varrho}}$  the balls with respect to the two metrics. Define  $V(z, w) = |B(z, \varrho(z, w))|$ ,  $\tilde{V}(z, w) = |B(z, \tilde{\varrho}(z, w))|$ . Let  $K(z, w)$  denote the distribution-kernel for  $\mathcal{L}^{-1}$ , and  $S$  the Szegö kernel. Here  $z, w \in \mathbb{C}^2 \times \mathbb{R}$ .

### Theorem 6.1 [M1].

- $|K(z, w)| \leq C \left( \tilde{\varrho}(z, w)^2 / \tilde{V}(z, w) \right) \cdot \log \left( 2 + \frac{\tilde{\varrho}(z, w)}{\varrho(z, w)} \right)$ .
- *The last estimate becomes false if the logarithmic factor is removed.*
- $|S(z, w)| \leq C \tilde{V}(z, w)^{-1}$ .
- $|L_2 \bar{L}_2 K(z, w)| \leq C V(z, w)^{-1}$ .

It is primarily  $\tilde{\varrho}$ , rather than  $\varrho$ , which figures in the first two estimates. Two further surprises: the extra singularity present in  $K$  disappears when appropriate derivatives are applied to it to yield  $S$ , while application of  $L_2 \bar{L}_2$  yields a kernel which satisfies natural estimates with respect to  $\varrho$  rather than  $\tilde{\varrho}$ .

Consider a pseudoconvex domain  $\Omega \subset \mathbb{C}^{n+1}$  of finite type, on which the Levi form has at most one degenerate eigenvalue at each boundary point. Fix smooth anti-holomorphic vector fields  $\bar{L}_1, \dots, \bar{L}_n$  which form a basis for the anti-holomorphic tangent space to  $\partial\Omega$  in some open set, with respect to which basis the Levi form is diagonalized, and such that the Levi form restricted to the span of  $\bar{L}_2, \dots, \bar{L}_n$  is nondegenerate. Let  $\tilde{\varrho}$  be the metric associated to the [FP]-type operator  $\tilde{\mathcal{L}} = -(X_1^2 + Y_1^2 + \sum_{j \geq 2} [X_j \lambda X_j + Y_j \lambda Y_j])$ , where  $\lambda$  denotes the degenerate eigenvalue. Define  $\tilde{V}$  in terms of  $\tilde{\varrho}$  as usual.

**Theorem 6.2** [M2]. *Let  $\Omega$  be pseudoconvex and of finite type, and suppose that its Levi form has at most one degenerate eigenvalue. Then the Szegö kernel on  $\partial\Omega$  satisfies*

$$|S(z, w)| \leq C \tilde{V}(z, w)^{-1} \quad \forall z \neq w \in \partial\Omega.$$

Moreover the appropriate bounds hold for derivatives, so that  $S$  is a singular integral operator with respect to the structure of a space of homogeneous type defined by  $\tilde{\varrho}$ . Thus it extends to an operator bounded on  $L^p$  for all  $p \in (1, \infty)$ ; it also improves Hölder classes.

It is easy to see via microlocal analysis why two different geometric structures ought to come into play. Let  $(\zeta, \tau) \in \mathbb{R}^4 \times \mathbb{R}$  be dual variables. Then microlocally where  $\tau > 0$  and  $|\tau| \geq C|\zeta|$ , one has for any  $f$ :

$$\begin{aligned} \langle \mathcal{L}P^+f, P^+f \rangle &= \|L_1P^+f\|^2 + \|\bar{L}_2P^+f\|^2 \\ &\sim \|L_1P^+f\|^2 + \|\bar{L}_1P^+f\|^2 + \langle \lambda|T|P^+f, P^+f \rangle + \|\bar{L}_2P^+f\|^2 \end{aligned}$$

modulo lower-order terms. One might also hope to estimate  $L_2P^+f$ , but

$$\|\bar{L}_2P^+f\|^2 - \|L_2P^+f\|^2 = \langle [\bar{L}_2, L_2]P^+f, P^+f \rangle \sim -\langle |T|P^+f, P^+f \rangle$$

(modulo a constant factor and zero-order terms) is potentially large and negative, and is not controlled by  $\langle \mathcal{L}P^+f, P^+f \rangle$ . The best that one can control is  $\langle \lambda[\bar{L}_2, L_2]P^+f, P^+f \rangle$ , so that

$$\begin{aligned} \langle \mathcal{L}P^+f, P^+f \rangle &\sim \|L_1P^+f\|^2 + \|\bar{L}_1P^+f\|^2 + \|\sqrt{\lambda}L_2P^+f\|^2 + \|\bar{L}_2P^+f\|^2 \\ &\sim \langle \tilde{\mathcal{L}}P^+f, P^+f \rangle + \|\bar{L}_2P^+f\|^2, \end{aligned}$$

again modulo lower-order terms. Thus  $\tilde{\mathcal{L}}$  enters the picture. On the other hand, in any compact subset of  $\mathbb{C}^2 \times \mathbb{R}$ ,  $\lambda$  is bounded above and one obtains

$$\langle \mathcal{L}P^-f, P^-f \rangle \sim \langle -(X_1^2 + Y_1^2 + X_2^2 + Y_2^2)P^-f, P^-f \rangle.$$

At first the estimates of Theorem 6.2 for  $S$  might appear improbable, given the extra singularity exhibited by  $K$ . However, the Szegö projection is equal to the orthogonal projection onto the intersection of the kernels of  $\bar{L}_1$  and  $\sqrt{\lambda}\bar{L}_2$ , and is thus heuristically expressible also as  $I - \tilde{\partial}_b^* \tilde{\square}_b^{-1} \tilde{\partial}_b$  where  $\tilde{\partial}_b f = (\bar{L}_1 f) \bar{\omega}_1 + \sqrt{\lambda}(\bar{L}_2 f) \bar{\omega}_2$  and  $\tilde{\partial}_b^*, \tilde{\square}_b^{-1}$  are the associated operators. (The same holds with  $\tilde{\partial}_b$  defined to be  $(\bar{L}_1 f) \bar{\omega}_1 + \gamma(\bar{L}_2 f) \bar{\omega}_2$ , for any constant  $\gamma \in \mathbb{R}$ .) Thus the Szegö kernel can be described entirely in the “ $\tilde{\varrho}$  – category”, and Theorem 6.2 becomes less surprising; this is part of the idea in [M2].

Concerning Hölder estimates one has the following:

**Theorem 6.3** [FKM]. *Let  $\Omega$  be smoothly bounded, pseudoconvex and of finite type  $m$ . Assume that its Levi form is (smoothly) diagonalizable in a neighborhood of each boundary point. Then for any  $0 < \alpha, \beta \in \mathbb{R} \setminus \mathbb{Z}$  satisfying  $\beta < \alpha + 2m^{-1}$ , for any  $(0, 1)$ -form  $f$  on  $\partial\Omega$ ,*

$$\|f\|_{A_\beta} \lesssim \|\square_b f\|_{A_\alpha} + \|f\|_{L^2}.$$

Moreover for  $\beta < 2/m$ ,  $\|f\|_{A_\beta} \lesssim \|\square_b f\|_{L^\infty}$ .

The proof is far too involved to be summarized here. Corresponding results hold for the  $\bar{\partial}$ -Neumann problem, and for the Szegö projection.<sup>7</sup>

## 7. Two Applications

### A. Characterization of Zero Varieties for the Nevanlinna Class

Let  $\Omega \subset \mathbb{C}^2$  be smooth, pseudoconvex and of finite type. Let  $r$  be a defining function,  $\Omega = \{r < 0\}$ . Let  $\Omega_\varepsilon = \{r < -\varepsilon\}$ . The Nevanlinna class is the set of all holomorphic functions on  $\Omega$  satisfying  $\sup_\varepsilon \int_{\partial\Omega_\varepsilon} \log^+ |f| < \infty$ . The zero variety  $Z \subset \Omega$  of a holomorphic function is said to satisfy the Blaschke condition if  $\int_Z |r(z)| d\sigma(z) < \infty$ , where  $\sigma$  is the induced volume element on  $Z$ .

**Theorem 7.1** [CNS2]. *Let  $Z \subset \Omega$  be a complex subvariety of codimension one, which satisfies the Blaschke condition. Then  $Z$  is the zero variety of a function in the Nevanlinna class.*

The converse was previously known and follows from Green's theorem. In the strictly pseudoconvex case this was proved earlier [He1, He2, Sk].<sup>8</sup>

The proof proceeds by the method of [L], in which one solves  $i\partial\bar{\partial}\Theta = [Z]$  where  $[Z]$  denotes the  $(1, 1)$  current defined by integration over  $Z$ . The desired Nevanlinna class function  $f$  will satisfy  $\Theta = \log |f|$ , so that one seeks a solution  $\Theta \in L^1(\partial\Omega)$ . Control of the  $L^1$  norm at the boundary comes from:

**Theorem 7.2** [CNS2].  $\|\bar{\partial}^* Nf\|_{L^1(\partial\Omega)} \leq C \|f\|_{L^1(\Omega)} + \|\frac{\mu}{r} f \wedge \bar{\partial}r\|_{L^1(\Omega)}$ .

The invariant  $\mu \gg r$  is defined in (5.1).

In order to apply Theorem 7.2, one needs an estimate of [BC] to the effect that the Blaschke condition automatically implies a formally stronger estimate for  $[Z]$ , generalizing a result of [Mal] for the strictly pseudoconvex case. Moreover the higher is the type of a given boundary point, the *stronger* becomes this estimate near that point, in such a fashion that the second term on the right-hand side in the theorem is exactly under control. The full strength of the precise estimates for the Neumann operator is used; no loss of  $\varepsilon$  in regularity can be afforded.

### B. Embeddability of Compact Three-Dimensional CR Manifolds.

Throughout this section  $M$  denotes a smooth, compact three-dimensional CR manifold without boundary. A CR embedding of  $M$  in  $\mathbb{C}^n$  is a smooth embedding,

---

<sup>7</sup> It is not known to this author whether one may take  $\beta = \alpha + 2m^{-1}$ , nor whether similar  $L^p$  estimates hold.

<sup>8</sup> An additional hypothesis imposed on  $Z$  in [Sk] is satisfied in this case since  $H^2(\Omega, \mathbb{C}) = 0$ .

all of whose component functions are CR, that is, are annihilated by  $\bar{\partial}_b$ . There exist arbitrarily small, real analytic perturbations of the standard structure on  $S^3$  which are not globally embeddable [Ro]. There exist strictly pseudoconvex CR structures, with the property that every function which is CR in a neighborhood of the origin is constant there [Ni].

There has been some interesting recent work in the negative direction. [BE] have proved that generic small perturbations of the standard structure on  $S^3$  are non-embeddable, even in the real-analytic category.  $\square_b^0 = \bar{\partial}_b^* \bar{\partial}_b$  is a self-adjoint operator which has (in the embeddable case) 0 as an eigenvalue of infinite multiplicity; the set of all positive eigenvalues is a locally discrete subset of the open interval  $(0, \infty)$ . Closed range fails if and only if there is a sequence of nonzero eigenvalues tending to 0, and [BE] show that generic perturbations create such eigenvalues.

[F] has given a simple construction of non-embeddable structures. One finds a simply connected, strictly pseudoconvex two-dimensional complex manifold  $\Omega$ , whose boundary admits a nontrivial finite cover,  $M$ , and one pulls the CR structure from  $\partial\Omega$  back to  $M$ . If embedded,  $M$  must bound a complex manifold, and the key is to prove that such a manifold would then be a nontrivial cover of  $\Omega$ .<sup>9</sup>

A very simple construction of non-embeddable (global and local) examples is in [R]. Further pathologies are exhibited there.

A combination of results of [Bu] and [K3] yields embeddability provided  $\bar{\partial}_b$  has closed range in  $L^2$ . Closed range would be a consequence of a subelliptic estimate for  $\square_b^0$ ; but  $\square_b^0$  is not subelliptic in dimension three.<sup>10</sup> In view of Theorem 2.2, closed range is equivalent to embeddability for strictly pseudoconvex three-dimensional CR manifolds. This result extends to the case of finite type:

**Theorem 7.3** [Ch4]. *A compact, pseudoconvex three-dimensional manifold without boundary and of finite type is embeddable, provided  $\bar{\partial}_b$  has closed range in  $L^2$ .*

To prove this one must show that there exist sufficiently many CR functions to provide local coordinates, and to separate points. In the strictly pseudoconvex case the latter is proved as follows: given  $x \neq y \in M$ , one first constructs a one-parameter family  $\{g_t : 0 < t < 1\}$  of functions satisfying  $g_t(x) = 1$ ,  $g_t(y) = 0$ , and  $\|\bar{\partial}_b g_t\|_{C^N} \rightarrow 0$  for any  $N$ . This is easily done by a formal power series argument, which relies on the existence of holomorphic support functions. From [K3] it follows that for sufficiently large  $\ell$ , one can solve  $\bar{\partial}_b u_t = \bar{\partial}_b g_t$  with  $\|u_t\|_{C^\ell} \leq C \|\bar{\partial}_b g_t\|_{C^\ell}$ . Then  $f_t = g_t - u_t$  is CR, and  $f_t(x) \rightarrow 1$  while  $f_t(y) \rightarrow 0$ .

There exist [KN] real analytic pseudoconvex hypersurfaces of finite type in  $\mathbb{C}^2$  for which any local holomorphic support function at a point  $x$  must vanish

<sup>9</sup> There do not exist complex manifolds of dimension greater than two with the required property.

<sup>10</sup> It is in higher dimensions, given strict pseudoconvexity, so that higher-dimensional compact, strictly pseudoconvex CR manifolds without boundary are always embeddable [Bt].

to infinite order; such functions cannot be used to construct adequate  $g_t$ . To circumvent this the strategy was to:

- show that on model hypersurfaces  $\{\text{Im}(z_2) = P(z_1)\}$ , where  $P$  is any homogeneous subharmonic, nonharmonic real-valued polynomial, there exists a CR function  $\varphi$  in the Schwartz class such that  $\varphi(0) \neq 0$ ,
- show that at any point of a pseudoconvex CR manifold of finite type, there exists a close approximation by one of the model domains,
- pull back a one-parameter family of dilates of  $\varphi$  from the model to  $\{g_t\}$  on the given manifold, in such a way that  $g_t(x) = 1$ ,  $g_t(y) = 0$ , and  $\|\bar{\partial}_b g_t\|_{L^q} \rightarrow 0$  for some finite  $q$ , and
- solve  $\bar{\partial}_b u_t = \bar{\partial}_b g_t$  with  $\|u_t\|_{C^0} \leq C \|\bar{\partial}_b g_t\|_{L^q}$ .

Then  $g_t - u_t$  again does the job, for small  $t$ . In the last step one does not know how to make  $\|\bar{\partial}_b g_t\|_{C^N}$  small for large  $N$ , and the  $L^p$  regularity theory is required to complete the proof.<sup>11</sup>

## 8. Analytic Hypoellipticity

In analogy with the  $C^\infty$  theory we say that  $\bar{\partial}_b$  is (relatively) analytic hypoelliptic on  $M$  if whenever  $\bar{\partial}_b u$  is real analytic in an open set  $V$  and  $u = \bar{\partial}_b^* v$  for some  $v \in L^2$  in  $V$ ,  $u$  must be analytic in  $V$ .

**Theorem 8.1** [G]. *Let  $\Omega \Subset \mathbb{C}^2$  be strictly pseudoconvex with real analytic boundary. Then  $\bar{\partial}_b$  is (relatively) analytic hypoelliptic, and the Szegö kernel is analytic off the diagonal on  $\partial\Omega \times \partial\Omega$ .*

This issue is not currently well understood for domains of finite type, but there are some counterexamples: Let  $\Omega = \{\text{Im}(w) = [\text{Re}(z)]^m\}$  for  $m = 4, 6, 8, \dots$

**Theorem 8.2** [CG].  *$\bar{\partial}_b$  is not (relatively) analytic hypoelliptic on  $\partial\Omega$ . Moreover the Szegö kernel is not analytic off the diagonal, nor is the Bergman kernel, restricted to  $\partial\Omega \times \partial\Omega$  minus the diagonal.*

In fact one has only Gevrey regularity of order  $m$  for general  $u$  as above, and in particular for the Szegö kernel. The proof rests on an explicit, though not entirely transparent, formula for the Szegö kernel in [N].

In higher dimensions we retain the standard definition of analytic hypoellipticity. Let  $n > 2$  and adopt coordinates  $(z', \zeta, w) \in \mathbb{C}^{n-2} \times \mathbb{C} \times \mathbb{C}$ .

**Theorem 8.3** [CG].  *$\square_b$  is not analytic hypoelliptic on  $\{\text{Im}(w) = |z'|^2 + [\text{Re}(\zeta)]^m\}$  for  $m = 4, 6, 8, \dots$  Nor is the Szegö kernel analytic off the diagonal.*

---

<sup>11</sup> The existence of continuous peak functions may also be inferred from this argument; a stronger result was obtained some time ago in [BF], and another proof has been obtained in [FSi].

Various subelliptic second-order differential operators with analytic coefficients are known not to be analytic hypoelliptic, but these examples seem to be new. An equivalent formulation is that the second-order operator  $\bar{\partial}_b \circ \bar{\partial}_b^*$  fails to be analytic hypoelliptic, in the standard sense, microlocally in a cone in which it is  $C^\infty$  hypoelliptic.

**Conjecture 8.4** [CG]. *On a hypersurface  $\{\text{Im}(z_2) = P(z_1)\}$ , where  $P$  is a subharmonic, nonharmonic polynomial,  $\bar{\partial}_b$  is (relatively) analytic hypoelliptic if and only if all zeroes of  $\Delta P$  are isolated in  $\mathbb{C}^1$ .*

## 9. Weighted Estimates for $\bar{\partial}$

The same philosophy as for the main results concerning  $\bar{\partial}_b$  may be applied in other situations. Let  $\varphi : \mathbb{C} \mapsto \mathbb{R}$  be subharmonic. Let  $\mu = \Delta\varphi$ , a positive, locally finite measure. Assume that

- (1)  $\mu$  is a doubling measure:  $\mu(B(z, 2r)) \leq C\mu(B(z, r)) \quad \forall z \in \mathbb{C}, r > 0$ .
- (2) There exists  $\delta > 0$  such that  $\mu(B(z, 1)) \geq \delta > 0$  for all  $z$ .

(1) should be viewed as a finite type condition; (2) replaces the usual assumption that  $\Delta\varphi \geq \delta$ . Then for each  $f \in L^2(\mathbb{C}, e^{-2\varphi})$  there exists  $u \in L^2(\mathbb{C}, e^{-2\varphi})$  satisfying  $\bar{\partial}u = f$  [Ch4, Ch6]. Let  $R$  be the bounded linear operator which assigns to each  $f$  the solution  $u$  with minimal norm in  $L^2(\mathbb{C}, e^{-2\varphi})$ , and denote also by  $R$  its distribution-kernel.

Let  $\lambda : \mathbb{C} \mapsto \mathbb{R}^+$  be  $C^\infty$  and satisfy  $\mu(B(z, \lambda(z))) \sim 1$  for all  $z$ . Define a Riemannian metric by  $d\varrho^2 = \lambda^{-2} ds^2$  where  $ds^2$  is the Euclidean metric.

**Theorem 9.1** [Ch6]. *Assume that  $\varphi$  satisfies (1) and (2). Then there exist  $C < \infty$ ,  $\varepsilon > 0$  such that for all  $z \neq \zeta \in \mathbb{C}$ ,*

$$|R(z, \zeta)| \leq C|z - \zeta|^{-1} e^{-\varepsilon\varrho(z, \zeta)} e^{\varphi(\zeta) - \varphi(z)}.$$

Moreover  $C, \varepsilon$  depend only on the constant in (1).

**Corollary 9.2** [Ch6]. *Let  $\varphi$  satisfy (1) and (2). Let  $p \in [1, \infty]$  and assume that  $f \in L^p(\mathbb{C}, e^{-\varphi})$ . Then there exists  $u$  satisfying  $\bar{\partial}u = f$ , such that  $u, \lambda^{-1}u \in L^p(\mathbb{C}, e^{-\varphi})$ .*

These results are closely related to the analysis of  $\bar{\partial}_b$  on the model hypersurfaces  $\{\text{Im}(z_2) = P(z_1)\}$ ; if  $P$  is a subharmonic but nonharmonic polynomial, then  $\varphi = P$  satisfies our hypotheses. Taking a partial Fourier transform in  $\text{Re}(z_2)$  reduces  $\bar{\partial}_b$  to  $\bar{\partial} + \tau \cdot P_z$  on  $\mathbb{C}^1$ , where  $\tau$  is a real parameter; this last equals  $e^{-\tau p} \circ \bar{\partial} \circ e^{\tau p}$ , so that studying it on  $L^2(\mathbb{C})$  with respect to Lebesgue measure is equivalent to studying the ordinary  $\bar{\partial}$  operator with a weight. The counterexamples of the last section become entirely natural from this perspective.

$L^p$  solvability of  $\bar{\partial}$  on bounded domains in  $\mathbb{C}^1$ , with arbitrary subharmonic  $\varphi$ , has been studied in [FiS2, Ber] and [A].

## References

- [A] E. Amar: Manuscript
- [BF] E. Bedford, J. E. Fornæss: A construction of peak functions on weakly pseudoconvex domains. *Ann. Math.* **107** (1978) 555–568
- [Be] J. Belanger: Hölder estimates for  $\bar{\partial}$ . 1987 Ph.D. thesis, Princeton University
- [Ber] B. Berndtsson: Weighted estimates for  $\bar{\partial}$  in domains in  $\mathbb{C}^n$ . *Duke Math. J.* (to appear)
- [BS] H. Boas, M. C. Shaw: Sobolev estimates for the Lewy operator on weakly pseudoconvex boundaries. *Math. Ann.* **274** (1986) 221–231
- [BC] A. Bonami, P. Charpentier: Estimations des  $(1,1)$  courants positifs fermés dans les domaines de  $\mathbb{C}^n$ . In: *Lecture Notes in Mathematics*, vol. 1094. Springer, Berlin Heidelberg New York 1984, pp. 44–52
- [BL] A. Bonami, N. Lohoué: Projecteurs de Bergman et Szegö pour une classe de domaines faiblement pseudoconvexes et estimations  $L^p$ . *Comp. Math.* **46** (1982) 159–226
- [Bt] L. Boutet de Monvel: Intégration des équations de Cauchy-Riemann induites formelles. Séminaire Goulaouic-Lions-Schwartz, Exposé IX
- [BSj] L. Boutet de Monvel, J. Sjöstrand: Sur la singularité des noyaux de Bergman et de Szegö. *Asterisque* **34–5** (1976) 123–164
- [Bu] D. Burns: Global behavior of some tangential Cauchy-Riemann equations. PDE and Geometry Conference, Park City, Utah, 1977. Dekker, New York 1979, pp. 51–56
- [BE] D. Burns, C. Epstein: Embeddability for three dimensional CR-manifolds. Preprint
- [C1] D. W. Catlin: Necessary conditions for subellipticity of the  $\bar{\partial}_b$ -Neumann problem. *Ann. Math.* **117** (1983) 147–171
- [C2] D. W. Catlin: Subelliptic estimates for the  $\bar{\partial}$ -Neumann problem on pseudoconvex domains. *Ann. Math.* **126** (1987) 131–191
- [C3] D. W. Catlin: Estimates of invariant metrics on pseudoconvex domains of domain two. *Math. Z.* **200** (1989) 429–466
- [C4] D. W. Catlin: A Newlander-Nirenberg theorem for manifolds with boundary. *Mich. Math. J.* **35** (1988) 233–240
- [CNS1] D. C. Chang, A. Nagel, E. M. Stein: Estimates for the  $\bar{\partial}$ -Neumann problem for pseudoconvex domains in  $\mathbb{C}^2$  of finite type. *Proc. Nat. Acad. Sci. USA* **85** (1988) 8771–8774
- [CNS2] D. C. Chang, A. Nagel, E. M. Stein: Estimates for the  $\bar{\partial}$ -Neumann problem in pseudoconvex domains of finite type in  $\mathbb{C}^2$ . Preprint
- [Ch1] M. Christ: Regularity properties of the  $\bar{\partial}_b$  equation on three-dimensional CR manifolds. *J. AMS* **1** (1988) 587–646
- [Ch2] M. Christ: Pointwise bounds for the relative fundamental solution of  $\bar{\partial}_b$ . *Proc. AMS* **104** (1988) 787–792
- [Ch3] M. Christ: Estimates for the  $\bar{\partial}_b$  equation and Szegö projection on CR manifolds. In: *Lecture Notes in Mathematics*, vol. 1384, pp. 146–158
- [Ch4] M. Christ: Embedding compact three-dimensional CR manifolds of finite type in  $\mathbb{C}^n$ . *Ann. Math.* **129** (1989) 195–213
- [Ch5] M. Christ: On the  $\bar{\partial}_b$  equation on three-dimensional CR manifolds. *Proc. Symp. Pure Math.* (to appear)
- [Ch6] M. Christ: On the  $\bar{\partial}$  equation in  $\mathbb{C}^1$  with weights. *J. Geom. Anal.* (to appear)
- [Ch7] M. Christ: Lectures on Singular Integral Operators. NSF-CBMS regional conference series, vol. 77. AMS, Providence 1990

- [Ch8] M. Christ: Estimates for fundamental solutions of second-order subelliptic differential operators. Proc. AMS **105** (1989) 166–172
- [CG] M. Christ, D. Geller: Counterexamples for analytic hypoellipticity on domains of finite type. Ann. Math. (to appear)
- [CW] R. R. Coifman, G. Weiss: Analyse Harmonique Non-Commutative Sur Certains Espaces Homogènes. (Lecture Notes in Mathematics, vol. 242). Springer, Berlin Heidelberg New York 1971
- [D] J. D'Angelo: Real hypersurfaces, orders of contact, and applications. Ann. Math. **115** (1982) 615–637
- [DT] M. Derridj, D. S. Tartakoff: Local analyticity for  $\square_b$  and the  $\bar{\partial}$ -Neumann problem at certain weakly pseudoconvex points. Comm. PDE **13** (1988) 1521–1600
- [DJ] G. David, J.-L. Journé: A boundedness criterion for generalized Calderón-Zygmund operators. Ann. Math. **120** (1984) 371–397
- [F] E. Falbel: Non-embeddable CR-manifolds and surface singularities. Preprint
- [Fe1] C. Fefferman: The Bergman kernel and biholomorphic mappings of pseudoconvex domains. Invent. math. **26** (1974) 1–66
- [Fe2] C. Fefferman: The uncertainty principle. Bull. AMS **9** (1983)
- [FK1] C. Fefferman, J. J. Kohn: Hölder estimates on domains of complex dimension two and on three dimensional CR manifolds. Adv. Math. **69** (1988) 223–303
- [FK2] C. Fefferman, J. J. Kohn: Estimates of kernels on three dimensional CR manifolds. Rev. Mat. Iberoam. **4** (1988) 355–405
- [FKM] C. Fefferman, J. J. Kohn, M. Machedon: Hölder estimates on CR manifolds with a diagonalizable Levi form. Rev. Mat. Iberoam. **4** (1988) 1–90
- [FP] C. Fefferman, D. H. Phong: Subelliptic eigenvalue problems. In: Proceedings, Conference on Harmonic Analysis in honor of Antoni Zygmund. Wadsworth, Belmont, CA, 1981, pp. 590–606
- [FSa] C. Fefferman, A. Sanchez-Callé: Fundamental solutions for second order subelliptic operators. Ann. Math. **124** (1986) 247–272
- [FS] G. B. Folland, E. M. Stein: Estimates for the  $\bar{\partial}_b$  complex and analysis on the Heisenberg group. Comm. Pure Appl. Math. **27** (1974) 429–522
- [Fo] J. Fornæss: Sup-norm estimates for  $\bar{\partial}$  in  $\mathbb{C}^2$ . Ann. Math. **123** (1986) 335–345
- [FSi1] J. Fornæss, N. Sibony: Construction of p.s.h. functions on weakly pseudoconvex domains. Duke Math. J. **59** (1989) 633–655
- [FSi2] J. Fornæss, N. Sibony: On  $L^p$  estimates for  $\bar{\partial}$ . Preprint, Université de Paris-Sud, 1989
- [G] D. Geller: Analytic Pseudodifferential Operators for the Heisenberg Group and Local Solvability. Princeton University Press, Princeton, NJ 1990
- [GS1] P. Greiner, E. M. Stein: On the solvability of some differential operators of type  $\square_b$ . In: Several Complex Variables, Proceedings of International Conferences, Cortona, Italy 1976–7. Scuola Normale Superiore, Pisa
- [GS2] P. Greiner, E. M. Stein: Estimates for the  $\bar{\partial}$ -Neumann' problem. Princeton University Press, Princeton, NJ 1977
- [HN] B. Helffer, J. Nourrigat: Hypoellipticité Maximale pour des Opérateurs Polynômes de Champs de Vecteurs. Birkhäuser, Boston 1985
- [He1] G. M. Henkin: Lewy's equation and analysis on pseudoconvex manifolds. I. Russian Math. Surveys **32** (1977) 59–130
- [He2] G. M. Henkin: Lewy's equation and analysis on pseudoconvex manifolds. II. Math. USSR Sb. **102** (1977) 63–64
- [H] L. Hörmander: Hypoelliptic second order differential equations. Acta Math. **119** (1967) 147–171

- [K1] J. J. Kohn: Boundary behavior of  $\bar{\partial}$  on weakly pseudoconvex manifolds of dimension two. *J. Diff. Geom.* **2** (1972) 523–542
- [K2] J. J. Kohn: The range of the tangential Cauchy-Riemann operator. *Duke Math. J.* **53** (1986) 525–545
- [K3] J. J. Kohn: Estimates for  $\bar{\partial}_b$  on pseudoconvex CR manifolds. *Proc. Symp. Pure Math.* **43** (1985) 207–217
- [KN] J. J. Kohn, L. Nirenberg: A pseudoconvex domain not admitting a holomorphic support function. *Math. Ann.* **201** (1973) 265–268
- [L] P. Lelong: Fonctionnelles analytiques et fonctions entières. Les Presses de l’Université de Montréal, Montréal, Canada 1968
- [M1] M. Machedon: Estimates for the parametrix of the Kohn Laplacian on certain domains. *Invent. math.* **91** (1988) 339–364
- [M2] M. Machedon: Szegő kernels on pseudoconvex domains with one degenerate eigenvalue. *Ann. Math.* **128** (1988) 619–640
- [Ma] P. Malliavin: Fonctions de Green d’un ouvert strictement pseudoconvexe et de classe de Nevanlinna. *C. R. Acad. Sci. Paris* **278** (1974) 141–144
- [Mc] J. D. McNeal: Boundary behavior of the Bergman kernel in  $\mathbb{C}^2$ . *Duke Math. J.* **58** (1989) 499–512
- [N] A. Nagel: Vector fields and nonisotropic metrics. In: *Beijing Lectures in Harmonic Analysis*. Annals of Mathematics Studies 112. Princeton University Press, Princeton, NJ 1986, pp. 241–306
- [NRSW1] A. Nagel, J.-P. Rosay, E. M. Stein, S. Wainger: Estimates for the Bergman and Szegő kernels in certain weakly pseudoconvex domains. *Bull. AMS* **18** (1988) 55–59
- [NRSW2] A. Nagel, J.-P. Rosay, E. M. Stein, S. Wainger: Estimates for the Bergman and Szegő kernels in  $\mathbb{C}^2$ . *Ann. Math.* **129** (1989) 113–149
- [NSW] A. Nagel, E. M. Stein, S. Wainger: Balls and metrics defined by vector fields, I: Basic properties. *Acta Math.* **155** (1985) 103–147
- [Ni] L. Nirenberg: Lectures on linear partial differential equations. CBMS-NSF regional conference series, vol. 17, AMS 1973
- [Ra] M. Range: Integral kernels and Hölder estimates for  $\bar{\partial}$  on pseudoconvex domains of finite type in  $\mathbb{C}^2$ . Preprint
- [R] J. P. Rosay: New examples of non-locally embeddable CR structures (with no non-constant CR distributions). *Annales Institut Fourier* **39** (1989) 811–823
- [Ro] H. Rossi: Attaching analytic spaces to a space along a pseudo-convex boundary. In: *Proceedings, Conference on Complex Analysis*. Springer, New York Berlin 1965, pp. 242–256
- [RS] L. P. Rothschild, E. M. Stein: Hypoelliptic differential operators and nilpotent groups. *Acta Math.* **137** (1976) 247–320
- [Sa] A. Sanchez-Calle: Fundamental solutions and geometry of the sum of squares of vector fields. *Invent. math.* **78** (1984)
- [S] M.-C. Shaw: Optimal Hölder and  $L^p$  estimates for  $\bar{\partial}_b$  on the boundaries of real ellipsoids in  $\mathbb{C}^n$ . *Trans. AMS* **324** (1991) 213–234
- [Si1] N. Sibony: Un exemple de domaine pseudoconvexe régulier où l’équation  $\bar{\partial}u = f$  n’admet pas de solution bornée pour  $f$  bornée. *Invent. math.* **62** (1980) 235–242
- [Si2] N. Sibony: On Hölder estimates for  $\bar{\partial}$ . Preprint, Université de Paris-Sud 1988
- [Sk] H. Skoda: Valeurs au bord pour les solutions de l’opérateur  $d''$ , et caractérisation des zeros des fonctions de la classe de Nevanlinna. *Bull. Soc. Math. France* **104** (1976) 225–299

- [Sm] H. Smith: A calculus for three-dimnnsional CR manifolds of finite type.  
Preprint
- [St] E. M. Stein: Singular Integrals and Differentiability Properties of Functions.  
Princeton University Press, Princeton, NJ 1970
- [W] S. Webster: On the proof of Kuranishi's embedding theorem. Ann. Inst. H. Poincaré Anal. Non Linéaire **6** (1989) 183–207



# Adapted Multiresolution Analysis, Computation, Signal Processing and Operator Theory

Ronald R. Coifman

Department of Mathematics, Yale University, New Haven, CT 06520, USA

We would like to describe recent developments, relating signal processing, numerical analysis and harmonic analysis.

Some of the main tasks of these fields involve efficient description of various classes of functions or transformations on functions. The usual problems encountered require handling of smoothness, oscillation or scaling patterns, and obtaining efficient representations (or coding) in a small number of parameters.

As for transformations on functions, some of the main problems concerning harmonic analysts involve preservation of various function spaces, spectral theory and operator calculus. The numerical analyst is mostly concerned with efficient algorithms for computing the effect of an operator, or of its inverse.

Both activities require a detailed understanding of the action of an operator on functions. In order to prove estimates or compute efficiently a good knowledge of the geometry and cancellation effects are necessary.

We will start by describing various methods for efficient analysis or description of functions. In the second part this analysis will be used to obtain fast numerical algorithms linking ideas and problems from classical analysis, such as Littlewood-Paley theory, to the ability to compute.

In conclusion, it will become evident that efficient compression methods, i.e. methods in which an operator or functions are described by as small a number of parameters as is possible (for a given precision), are directly related to our ability to compute fast, as well as to our analytic understanding. Most of the results discussed in this talk were obtained in collaboration with Yves Meyer, Vladimir Rokhlin, Gregory Beylkin, and Victor Wickerhauser.

## § 1. Data Compression, Orthonormal Bases, and Best Basis

We describe various constructions of orthonormal bases in function spaces (or sample spaces) which will permit an adaptation of bases to given functions or signals.

Over the last fifty years we have seen within mathematics, the need to modify the Fourier transform by permitting regions or bands of frequencies to be lumped together. The most common analysis of this type, the so-called Littlewood-Paley theory, (as developed by Marcinkiewicz, Zygmund, Calderon, Stein and many others) in which frequencies are grouped in dyadic intervals, has proved to be a

powerful and flexible tool. This method permits a blend of analysis in space and frequency simultaneously.

The recent discoveries by Stromberg [S], Meyer [M], Mallat and Daubechies [D] of various orthonormal wavelet bases has opened the door for using these methods (Littlewood-Paley analysis) in signal processing and numerical computation and, stimulated a range of discoveries and constructions of other classes of orthonormal bases. We refer the reader to Meyer's talk in these proceedings for further detail.

**Definitions of Modulated Waveform Libraries.** We now introduce the concept of a "Library of orthonormal bases". For the sake of exposition we restrict our attention to two classes of numerically useful waveforms, introduced recently by Y. Meyer and the author.

We start with trigonometric waveform libraries. These are localized sine transforms associated to covering by intervals of  $\mathbf{R}$  (more generally, of a manifold).

We consider a cover  $\mathbf{R} = \bigcup_{-\infty}^{\infty} I_i$ ,  $I_i = [\alpha_i, \alpha_{i+1}]$   $\alpha_i < \alpha_{i+1}$ , write  $\ell_i = \alpha_{i+1} - \alpha_i = |I_i|$  and let  $p_i(x)$  be a window function supported in  $[\alpha_i - \ell_{i-1}/2, \alpha_{i+1} + \ell_{i+1}/2]$  such that

$$\sum_{-\infty}^{\infty} p_i^2(x) = 1$$

and

$$p_i^2(x) = 1 - p_i^2(2\alpha_{i+1} - x) \quad \text{for } x \text{ near } \alpha_{i+1}$$

then the functions

$$S_{i,k}(x) = \frac{2}{\sqrt{2\ell_i}} p_i(x) \sin \left[ (2k+1) \frac{\pi}{2\ell_i} (x - \alpha_i) \right]$$

form an orthonormal basis of  $L^2(\mathbf{R})$  subordinate to the partition  $p_i$ . The collection of such bases forms a library of orthonormal bases.

It is easy to check that if  $H_{I_i}$  denotes the space of functions spanned by  $S_{i,k}$   $k = 0, 1, 2, \dots$  then  $H_{I_i} + H_{I_{i+1}}$  is spanned by the functions

$$P(x) \frac{1}{\sqrt{2(\ell_i + \ell_{i+1})}} \sin \left[ (2k+1) \frac{\pi}{2(\ell_i + \ell_{i+1})} (x - \alpha_i) \right]$$

where

$$P = (p_i^2(x) + p_{i+1}^2(x))^{1/2}$$

is a "window" function covering the interval  $I_i \cup I_{i+1}$ .

Another new library of orthonormal bases called the Wavelet packet library can be constructed. This collection of modulated wave forms, corresponds roughly to a covering of "frequency" space. This library contains the wavelet basis, Walsh functions, and smooth versions of Walsh functions called wavelet packets.

We'll use the notation and terminology of [D], whose results we shall assume.

We are given an exact quadrature mirror filter  $h(n)$  satisfying the conditions of Theorem (3.6) in [D], p. 964, i.e.

$$\sum_n h(n-2k)h(n-2\ell) = \delta_{k,\ell}, \quad \sum_n h(n) = \sqrt{2}.$$

We let  $g_k = h_{l-k}(-1)^k$  and define the operations  $F_i$  on  $\ell^2(\mathbf{Z})$  into “ $\ell^2(2\mathbf{Z})$ ”

$$(1.0) \quad \begin{aligned} F_0\{s_k\}(i) &= 2 \sum s_k h_{k-2i} \\ F_1\{s_k\}(i) &= 2 \sum s_k g_{k-2i}. \end{aligned}$$

The map  $\mathbf{F}(\mathbf{s}_k) = \mathbf{F}_0(\mathbf{s}_k) \oplus \mathbf{F}_1(\mathbf{s}_k) \in \ell^2(2\mathbf{Z}) \oplus \ell^2(2\mathbf{Z})$  is orthogonal and

$$(1.1) \quad F_0^* F_0 + F_1^* F_1 = I.$$

We now define the following sequence of functions.

$$(1.2) \quad \begin{cases} W_{2n}(x) = \sqrt{2} \sum h_k W_n(2x - k) \\ W_{2n+1}(x) = \sqrt{2} \sum g_k W_n(2x - k). \end{cases}$$

Clearly the function  $W_0(x)$  can be identified with the scaling function  $\varphi$  in [D] and  $W_1$  with the basic wavelet  $\psi$ .

Let us define  $m_0(\xi) = \frac{1}{\sqrt{2}} \sum h_k e^{-ik\xi}$  and

$$m_1(\xi) = -e^{i\xi} \bar{m}_0(\xi + \pi) = \frac{1}{\sqrt{2}} \sum g_k e^{ik\xi}.$$

*Remark.* The quadrature mirror condition on the operation  $\mathbf{F} = (\mathbf{F}_0, \mathbf{F}_1)$  is equivalent to the unitarity of the matrix

$$\mathcal{M} = \begin{bmatrix} m_0(\xi) & m_1(\xi) \\ m_0(\xi + \pi) & m_1(\xi + \pi) \end{bmatrix}.$$

Taking the Fourier transform of (1.2) when  $n = 0$  we get

$$\hat{W}_0(\xi) = m_0(\xi/2) \hat{W}_0(\xi/2)$$

i.e.,

$$\hat{W}_0(\xi) = \prod_{j=1}^{\infty} m_0(\xi/2^j)$$

and

$$\hat{W}_1(\xi) = m_1(\xi/2) \hat{W}_0(\xi/2) = m_1(\xi/2) m_0(\xi/4) m_0(\xi/2^3) \cdots$$

More generally, the relations (1.2) are equivalent to

$$(1.3) \quad \hat{W}_n(\xi) = \prod_{j=1}^{\infty} m_{e_j}(\xi/2^j)$$

and  $n = \sum_{j=1}^{\infty} e_j 2^{j-1}$  ( $e_j = 0$  or  $1$ ).

The functions  $W_n(x - k)$  form an orthonormal basis of  $L^2(\mathbf{R}^n)$ .

We define a *library* of wavelet packets to be the collection of functions of the form  $W_n(2^\ell x - k)$  where  $\ell, k \in \mathbf{Z}, n \in \mathbf{N}$ . Here, each element of the library is determined by a scaling parameter  $\ell$ , a localization parameter  $k$  and an oscillation

parameter  $n$ . (The function  $W_n(2^\ell x - k)$  is roughly centered at  $2^{-\ell}k$ , has support of size  $\approx 2^{-\ell}$  and oscillates  $\approx n$  times).

We have the following simple characterization of subsets forming orthonormal bases.

**Proposition.** *Any collection of indices  $(\ell, n)$  such that the intervals  $[2^\ell n, 2^\ell n + 1]$  form a disjoint cover of  $[0, \infty)$  gives rise to an orthonormal basis of  $L^2$ .*

Motivated by ideas from signal processing and communication theory V. Wickerhauser and the author were led to measure the “distance” between a basis and a function in terms of the Shannon entropy of the expansion. More generally, let  $H$  be a Hilbert space.

Let  $v \in H$ ,  $\|v\| = 1$  and assume

$$H = \bigoplus H_i$$

an orthogonal direct sum. We define

$$\varepsilon^2(v, \{H_i\}) = - \sum \|v_i\|^2 \ell n \|v_i\|^2$$

as a measure of distance between  $v$  and the orthogonal decomposition.

$\varepsilon^2$  is characterized by the Shannon equation which is a version of Pythagoras' theorem.

Let

$$H = \bigoplus (\sum H^i) \oplus (\sum H_j) = H_+ \oplus H_-$$

$H^i$  and  $H_j$  give orthogonal decomposition  $H_+ = \sum H^i, H_- = \sum H_j$ . Then

$$\varepsilon^2(v; \{H^i, H_j\}) = \varepsilon^2(v, \{H_+, H_-\}) + \|v_+\|^2 \varepsilon^2 \left( \frac{v_+}{\|v_+\|}, \{H^i\} \right) + \|v_-\|^2 \varepsilon^2 \left( \frac{v_-}{\|v_-\|}, \{H_j\} \right).$$

This is Shannon's equation for entropy (if we interpret as in quantum mechanics  $\|P_{H_+} v\|^2$  as the “probability” of  $v$  to be in the subspace  $H_+$ ).

This equation enables us to search for a smallest entropy space decomposition of a given vector.

In fact, for the example of the first library restricted to covering by dyadic intervals we can start by calculating the entropy of an expansion relative to a local trigonometric basis for intervals of length one, then compare the entropy of an adjacent pair of intervals to the entropy of an expansion on their union. Pick the expansion of minimal entropy and continue until a minimum entropy expansion is achieved.

In practice, discrete versions of this scheme can be implemented in  $CN \log N$  computations (where  $N$  is the number of discrete samples  $N = 2^L$ .)

For voice signals and images this procedure leads to remarkable compression algorithms (see [CMQW]).

Of course, while entropy is a good measure of concentration of an expansion, various other information cost functions are possible, permitting discrimination and choice between various special function expansion.

Other possible libraries can be constructed. The space of frequencies can be decomposed into pairs of symmetric windows around the origin, on which

a smooth partition of unity is constructed. This and other constructions were obtained by one of our students E. Laeng [L].

Higher dimensional libraries can also be easily constructed (as well as libraries on manifolds), leading to new and direct analysis methods for linear transformations.

## § 2. Wavelets, Wavelet Packets and Numerical Algorithms

The usual way to analyze an integral operator relies on the Fourier transform. In principle, for convolution operators the problem is solved, although as we well know it is essentially impossible to rely solely on the Fourier transform to understand the effect of the operator on various classes of functions (such as Hölder,  $L^p$  etc.).

For operators which are not of convolution type even the problem of proving  $L^2$  boundedness (for example), can be extremely difficult. We claim that good methods for fast computation of such operators shed light on their analysis and provide new approaches.

Initially the relations between computation and Calderón-Zygmund theory was pointed out to the author by V. Rokhlin who in his design of the Fast Multipole algorithm for computing potential interactions has essentially reinvented many of the ingredients of Calderón-Zygmund theory.

Rokhlin constructed a fast algorithm of order  $N$  to compute all sums

$$p_j = \sum_{i=1}^N \frac{g_i g_j}{|x_i - x_j|} \quad \text{where } x_i \in \mathbf{R}^3 \quad i = 1, \dots, N$$

although, naively it would seem to be impossible to do this calculation in less than  $N^2$  computations, since this is the number of interactions. He observed that the effect of a cloud of charges located in a box can be described to any accuracy by the effect of a single multipole at the center of the box, requiring only a few numbers (Taylor coefficients of the field at the center of external boxes removed from the source). He organized all boxes in a dyadic hierarchy enabling an efficient  $O(N)$  algorithm. (This algorithm is  $O(N)$  independently of the configuration of the charges providing therefore a substantial improvement over  $FFT$ ).

As we know, Littlewood-Paley and Calderón-Zygmund theory can achieve similar goals. Wavelet based algorithms providing an elegant reformulation and generalization of the multipole algorithms, were developed by Beylkin, Rokhlin and the author [BCR].

Before describing these methods in detail we return to the basic question of efficient computation of an integral operator.

$$Tf(x) = \int_0^1 k(x, y) f(y) dy \quad x \in [0, 1].$$

Clearly, if we can write

$$k(x, y) = \sum a_{\alpha\beta} w_\alpha(x) w_\beta(y)$$

where the number of  $a_{\alpha\beta}$  is small, we would reduce the computation to

$$\int f(y)w_\beta(y)dy = d_\beta$$

and to

$$\sum_\alpha \left( \sum_\beta a_{\alpha\beta} d_\beta \right) w_\alpha(x).$$

If we choose  $w_\beta$  to be the eigenvectors of  $T$  (say if  $T$  is self adjoint) at least the matrix  $a_{\alpha\beta}$  would be diagonal. The price of course would be to compute  $d_\beta$  and the sum efficiently (say like *FFT*). Such algorithms of course are not known in general. Instead, we can compromise. View  $k(x, y)$  as an image (i.e.  $k(x, y)$ ) would represent the light intensity at pixel  $(x, y)$ ) and try to compress the image. We are naturally led to consider a library such as the wavelet packets and to a selection of an orthonormal basis of  $L^2((0, 1) \times (0, 1))$  such that the matrix  $a_{\alpha\beta}$  has highest concentration (or lowest entropy).

Of course, when the kernel under consideration satisfies estimates invariant under translations and dilations. In particular, if  $k(x, y)$  is a Calderón-Zygmund operator or a pseudo differential operator we expect the best basis to have a similar behavior i.e., we are forced to consider the wavelet basis. Remarkably, this algorithm (see [BCR]) corresponds on the one hand to the *FMM* methods of Rokhlin [R] and, on the other hand, to the so-called  $P_t, Q_t$  analysis of Calderón-Zygmund operators and the  $T(1)$  theorem of David and Journé (see [DJ]).

Concretely these methods can be most simply described by using the Haar functions  $h_I(x)$  where  $h_I(x) = \frac{1}{|I|^{\frac{1}{2}}}$  on the left half of the dyadic interval  $I$  and  $h_I(x) = -\frac{1}{|I|^{\frac{1}{2}}}$  on the right half, and zero elsewhere. We also let

$$\chi_I = \frac{1}{|I|^{\frac{1}{2}}} \quad \text{on } I \quad \chi_I = 0 \quad x \notin I.$$

We expand  $k(x, y)$  in terms of the two dimensional Haar functions  $h_I(x)h_{I'}(y)$ ,  $h_I(x)\chi_{I'}(y)$ ,  $\chi_I(x)h_{I'}(y)$  as

$$k(x, y) = \sum \alpha_{II'} h_I(x)h_{I'}(y) + \sum \beta_{II'} h_I(x)\chi_{I'}(y) + \sum \gamma_{II'} \chi_I(x)h_{I'}(y)$$

where  $\alpha_{II'} = \iint k(x, y)h_I(x)h_{I'}(y)dxdy$

$$\beta_{II'} = \iint k(x, y)h_I(x)\chi_{I'}(y)dxdy$$

$$\gamma_{II'} = \iint k(x, y)\chi_I(x)h_{I'}(y)dxdy.$$

Introducing this representation of  $k$  we obtain

$$(1.3) \quad T(f)(x) = \sum_I h_I(x) \sum_{I'} \alpha_{II'} d_{I'} + \sum_I h_I(x) \sum_{I'} \beta_{II'} s_{I'} \sum_I \chi_I(x) \sum_{I'} \gamma_{II'} d_{I'}$$

where each sum in  $I'$  involves only dyadic intervals of length  $|I|$  and where  $s_I = \langle \chi_I, f \rangle$ ,  $d_I = \langle h_I, f \rangle$ . The matrix realisation of this computation (Fig. 1) does

not correspond to the ordinary realization of the operator in terms of the Haar basis. In fact the  $s_j$  are dependant on the coordinates  $d_j$ ; by doubling the size of the matrix we gain a block decomposition by scale (all interscale interactions occur through  $s_j$ ). When applied to specific classes of operators such as C-Z or pseudo-differential operators this procedure yields banded matrices, with band width depending on the desired precision and the choice of wavelet. If we let  $P_j$  denote the orthogonal projection on the space of functions constant on dyadic intervals of length  $2^{-j}$  and approximate  $T$  by  $T_n = P_n T P_n$ .

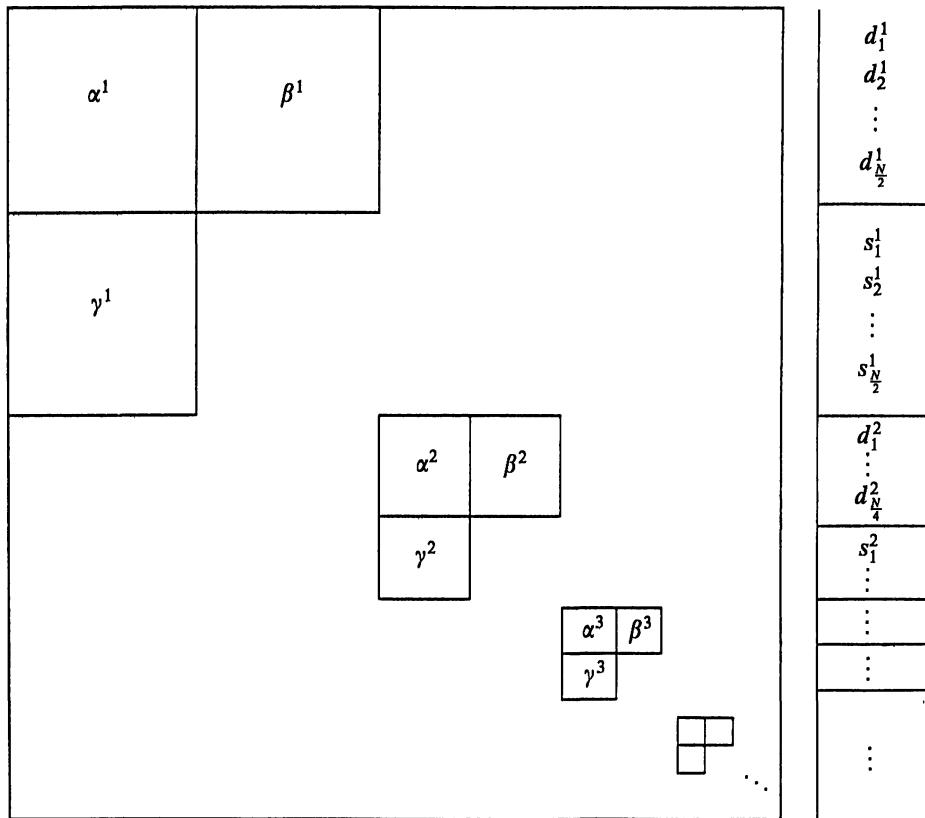


Fig. 1

We can rewrite

$$T_n = \sum_{j=1}^n (T_j - T_{j-1}) + T_0 = \sum_{j=1}^n (Q_j T Q_j + Q_j T P_j + P_j T Q_j) + T_0$$

where  $Q_j = P_j - P_{j-1}$ .

This decomposition relates the previous expansion to the Littlewood Paley approach and the proof of the "T of 1" theorem of David and Journe, see also [Se].

In general, we might hope for decompositions permitting us also to handle oscillatory integrals or higher singularities such as the one arising from wave propagation, Fourier integral operators, Radon transforms, and others.

It now seems clear that at least for some of these problems where previously microlocalization was the tool of choice for obtaining estimates, adapted orthonormal bases may become the tool for numerical computation.

Another class of examples relate to the Cauchy integral on chord arc curves. Here we study

$$C(f)(s) = p \cdot v \int \frac{f(t)}{z(s) - z(t)} dt$$

where  $s$  is the arclength parameter and it is assumed that

$$|z(s) - z(t)| > \delta |s - t| \quad \text{for some } \delta > 0.$$

Again, the problem of computing  $C(f)$  rapidly as described before seems at first impossible since the curve is not assumed to be more than once differentiable, and rapid decay reflects smoothness. Here however,  $z'(t)dt$  is a complex measure relative to which everything is smooth. Choosing an orthogonal basis relative to the complex measure  $dz$  with sufficiently many  $dz$  vanishing moments leads one to a rapid algorithm as well as to a simple, beautiful proof of the  $L^2(ds)$  boundedness of the Cauchy integral (see [CJS]). This fact is easily proved because the Cauchy operator becomes almost diagonal in this basis. In this problem we see a natural tie between the geometry of the curve and operator theory. Moreover the existence of a "good" basis in  $L^2(ds)$  is equivalent to the chord-arc condition.

Equivalent relations between the geometry of the curve and operator theory can be seen through the size of the wavelet expansion coefficients of  $z'(t)$ . These coefficients measure the deviation from flatness on various scales of the curve. Analogous methods, initiating nonlinear Littlewood-Paley theory, studying the deviation from flatness of general sets in  $\mathbb{R}^n$  have been introduced by P. Jones [PJ] in his beautiful characterization of subsets of rectifiable curves and by G. David and S. Semmes [DS] in their work on operator theory on surfaces.

## References

- [BCR] Beylkin, G., Coifman, R., Rokhlin, V.: Fast wavelet transforms and numerical algorithms I. To appear in Comm. Pure Appl. Math.
- [CJS] Coifman, R., Jones, P.W., Semmes, S.: Two elementary proofs of the  $L^2$  boundedness of Cauchy integrals on Lipschitz curves. J. AMS 2 (1989) 553–564
- [CMQW] Coifman, R., Meyer, Y., Quake, S., Wickerhauser, V.: Signal processing and compression with wavelet packets. To appear
- [D] Daubechies, I.: Orthonormal bases of compactly supported wavelets. Comm. Pure Appl. Math. XL1 (1988)
- [DJ] David, G., Journé, J.L.: A boundedness criterion for generalized Calderon-Zygmund operators. Ann. Math. 120 (2) (1984) 371–397
- [DS] David, G., Semmes, S.: Audela des graphies Lipschitziens. To appear in Asterisque-Au-delà
- [L] Laeng, E.: Une base orthonormale de  $L^2(\mathbb{R})$ .... To appear in C. R. Acad. Sci. Paris (1990)

- [M] Meyer, Y.: Principe d'incertitude, bases hilbertiennes et algébres d'opérateurs. Séminaire Bourbaki, 1985–85, 662. Astérisque (Société Mathématique de France)
- [PJ] Jones, P.: Rectifiable sets and the traveling salesman problem. Invent. math. (1990)
- [S] Stromberg, J.O.: A modified Haar system and higher order spline systems. Conference in Harmonic Analysis in Honor of Antoni Zygmund. Wadsworth Math. Series, ed. by W. Beckner et al., vol. II, pp. 475–493
- [Se] Semmes, S.: Nonlinear Fourier analysis. Bull. AMS **20** (1) (1989)



# Rational Maps and Kleinian Groups

Curt McMullen

Princeton University, Department of Mathematics, Fine Hall, Washington Road  
Princeton NJ 08544, USA

## 1. Introduction

There are many parallels between the theory of iterated rational maps  $f$  and that of Kleinian groups  $\Gamma$ , considered as dynamical systems on the Riemann sphere  $\widehat{\mathbb{C}}$ .<sup>1</sup> In this paper we will survey three chapters of this developing theory, and the Riemann surface techniques they employ:

1. The combinatorics of critically finite rational maps and the geometrization of Haken 3-manifolds via iteration on Teichmüller space.
2. Renormalization of quadratic polynomials and 3-manifolds which fiber over the circle.
3. Boundaries and laminations — Teichmüller space in Bers' embedding and the Mandelbrot set.

## 2. The Theme of Short Geodesics

What are the possible topological forms for a conformal dynamical system?

Part of the answer is provided by two theorems, due to Thurston, which employ *iteration on Teichmüller space* to construct rational maps and Kleinian groups of a given topological form. More precisely, the iteration either *finds a geometric model* or *reveals a topological obstruction* to its existence. This dichotomy stems from:

**Theorem 1** [Mum]. *Let  $X_n$  be a sequence of points in the moduli space  $\mathcal{M}_{g,k}$  of hyperbolic Riemann surfaces of genus  $g$  with  $k$  punctures. After passing to a subsequence, either*

- $X_n$  converges to  $X$  in  $\mathcal{M}_{g,k}$ , or
- there is a collection of disjoint simple closed geodesics  $S_n$  on  $X_n$  such that the hyperbolic length of  $S_n$  tends to zero.

---

<sup>1</sup> See [Sul2] for part of the dictionary.

## 2.1 Critically Finite Rational Maps

Let  $f : S^2 \rightarrow S^2$  be a branched covering of the sphere of degree greater than one, and let  $P$  denote the *post-critical set* of  $f$ , i.e.

$$P = \bigcup_{n=1}^{\infty} f^n(B),$$

where  $B$  denotes the branch points (at which  $f$  is locally many-to-one). If  $|P| < \infty$  we say  $f$  is *critically finite*. Two such maps  $f$  and  $g$  are *combinatorially equivalent* if there is a homeomorphism  $h : (S^2, P_f) \rightarrow (S^2, P_g)$  such that  $hfh^{-1}$  and  $g$  are isotopic rel  $P_g$ .

A critically finite map is a generalization to the complex domain of the *kneading sequence* for maps of the interval.

The following theorem provides a topological characterization of critically finite rational maps.

**Theorem 2** [Th3, DH3]. *Let  $f : S^2 \rightarrow S^2$  be critically finite with hyperbolic orbifold. Then either*

- *$f$  is combinatorially equivalent to a rational map  $g : \widehat{\mathbb{C}} \rightarrow \widehat{\mathbb{C}}$ , unique up to automorphisms of  $\widehat{\mathbb{C}}$ , or*
- *there is an  $f$ -invariant system of disjoint simple closed curves  $\Gamma$  in  $S^2 - P$  providing a topological obstruction to such an equivalence.*

The technical condition “with hyperbolic orbifold” rules out certain elementary cases (which are also understood). It is satisfied, for example, if  $|P| > 4$ .

*Sketch of the proof.* The space of Riemann surface structures on  $(S^2, P)$ , up to isotopy rel  $P$  is exactly the Teichmüller space of the sphere with  $|P|$  distinguished points, denoted  $\text{Teich}(S^2, P)$ . Given such a structure, pull it back by  $f$  to obtain a new structure on the same space: this defines a map

$$T_f : \text{Teich}(S^2, P) \rightarrow \text{Teich}(S^2, P).$$

A fixed point for  $T_f$  gives an invariant complex structure and therefore a rational map  $g$  combinatorially equivalent to  $f$ .

This iteration has two fundamental features:

- $T_f^k$  contracts the Teichmüller metric (for some fixed iterate  $k$ ); and
- the contraction at a point  $X$  in  $\text{Teich}(S^2, P)$  is less than  $c[X] < 1$  where  $c[X]$  is a continuous function depending only on the location of  $X$  in moduli space.

Now try to locate a fixed point of  $T_f$  by studying the sequence of iterates  $X_n = T_f^n(X_0)$  of an arbitrary starting guess  $X_0$ . If  $[X_n]$  returns infinitely often to a compact subset  $K$  of moduli space, then due to uniform contraction over  $K$ , the sequence converges to a fixed point and  $f$  is equivalent to a rational map.

Otherwise, by Mumford's theorem, the length of the shortest geodesic on  $X_n$  tends to zero. Set  $\Gamma = \{\text{isotopy classes of very short geodesics on } X_n\}$  for  $n$  sufficiently large. Since the Teichmüller distance from  $X_n$  to  $X_{n+1}$  is bounded, lengths change by only a bounded factor. Therefore  $\Gamma$  is  $f$ -invariant, in the sense that any geodesic representing a component of  $f^{-1}(\gamma)$  is again in  $\Gamma$ .

Let  $A : \mathbb{R}^\Gamma \rightarrow \mathbb{R}^\Gamma$  be defined by  $A_{\delta\gamma} = \sum_\alpha \deg(f : \alpha \rightarrow \gamma)^{-1}$ , where the sum is over components  $\alpha$  of  $f^{-1}(\gamma)$  homotopic to  $\delta$ . By analyzing the geometry of short geodesics, one shows the leading eigenvalue of  $A$  is  $\geq 1$ . This provides the desired topological obstruction.

Indeed, if  $f$  is equivalent to a rational map  $g$ , then one can thicken the curves in  $\Gamma$  to disjoint annuli, with conformal moduli  $m_\gamma > 0$ . By considering inverse images of these annuli under  $g$ , one finds the same curves can be represented by annuli with moduli  $m'_\delta > \sum_\gamma A_{\delta\gamma} m_\gamma$ . It follows that some curve can be thickened to an annulus of arbitrarily large conformal modulus, a contradiction.  $\square$

## 2.2 Haken 3-Manifolds

There is a parallel theory of iteration in Thurston's construction of hyperbolic structures on Haken manifolds. For simplicity we stick to closed manifolds.

To a Kleinian group  $\Gamma$  one associates the 3-dimensional *Kleinian manifold*  $N = (\mathbb{H}^3 \cup \Omega)/\Gamma$ , where  $\Omega \subset \widehat{\mathbb{C}}$  is the domain of discontinuity. Then  $N$  has a hyperbolic structure on its interior and a conformal structure on its boundary.

**Theorem 3** [Th1, Mor]. *Let  $M^3$  be a closed Haken 3-manifold. Then either*

- $M^3$  is diffeomorphic to a unique hyperbolic manifold  $\mathbb{H}^3/\Gamma$ , or
- there is a map of a torus into  $M^3$ , injective on  $\pi_1$ , providing a topological obstruction to a hyperbolic structure on  $M^3$ .

*Sketch of the proof.* We combine Thurston's original approach with the Riemann surface techniques of [Mc2] and emphasize the parallel with the geometrization of rational maps.

A 3-manifold is *Haken* if it can be constructed by starting with 3-balls, and repeatedly gluing along incompressible submanifolds of the boundary. The idea of the proof is to carry out the construction geometrically, at each stage providing the pieces with hyperbolic structures. An orbifold technique [Mor, Fig. 14.6] reduces the problem to the case of gluing along the entire boundary.

Iteration enters at the inductive step: given a compact 3-manifold  $M^3$  and gluing instructions encoded by an orientation-reversing involution  $\tau : \partial M^3 \rightarrow \partial M^3$ , we must construct a hyperbolic structure on  $M^3/\tau$ . By induction  $M^3$  is diffeomorphic to a Kleinian manifold  $N^3$ . Unlike the case of a closed manifold, which admits at most one hyperbolic structure by Mostow rigidity, the manifold  $N^3$  is flexible. The set of possible shapes for  $N^3$  is parameterized by the Teichmüller space of the boundary of  $M$ .

Which structure descends to  $M/\tau$ ? The answer can be formulated as a fixed point problem on Teichmüller space. Using the topology of  $M^3$ , Thurston defines

the *skinning map*

$$\sigma : \text{Teich}(\partial M) \rightarrow \text{Teich}(\overline{\partial M})$$

by forming quasifuchsian covering spaces for each component of the boundary, and recording the conformal structure on the new ends which appear. The gluing instructions determine an isometry

$$\tau : \text{Teich}(\overline{\partial M}) \rightarrow \text{Teich}(\partial M),$$

and a fixed point for

$$T = \sigma \circ \tau$$

solves the gluing problem.

Here the parallel with the construction of critically finite maps emerges. The completion of the proof will follow [Mc2].

Assume  $M^3$  is not an interval bundle over a surface (this special case is discussed in the next section). Then some fixed iterate  $T^k$  contracts the Teichmüller metric; in fact:

- The contraction of  $T^k$  at a point  $X$  in  $\text{Teich}(\partial M)$  is bounded by  $c[X] < 1$  where  $c[X]$  depends only on the location of  $X$  in moduli space.

As before, this reduces the proof to an analysis of short geodesics. Let  $X_n = T^n(X_0)$  be the forward orbit of a starting guess  $X_0$  in  $\text{Teich}(\partial M)$ . If  $[X_n]$  returns infinitely often to a compact subset of moduli space, the sequence converges and the gluing problem is solved.

Otherwise  $X_n$  develops short geodesics. With further analysis one finds these short geodesics bound cylinders in  $M^3$ , joined by  $\tau$  to form an incompressible torus in  $M^3/\tau$ . A closed hyperbolic manifold contains no such torus (it must correspond to a cusp), so we have located a topological obstruction to a hyperbolic structure.  $\square$

The bound on contraction  $c[X]$  comes from a general result in the theory Riemann surfaces.

Let  $Y \rightarrow X$  be a covering space of a hyperbolic Riemann surface  $X$  of finite area. Then there is a natural map  $\theta : \text{Teich}(X) \rightarrow \text{Teich}(Y)$ , defined by lifting complex structures from  $X$  to  $Y$ . Consider the case of the universal covering  $\mathbb{D} \rightarrow X$  where  $\mathbb{D}$  is the unit disk;  $G$  denotes the Fuchsian group of deck transformations.

**Theorem 4** [Mc1]. *The map  $\theta : \text{Teich}(X) \rightarrow \text{Teich}(\mathbb{D})$  is a contraction for the Teichmüller metric. Moreover  $\|d\theta\| < c[X] < 1$  where  $c$  depends continuously on the location of  $X$  in moduli space.*

This theorem is related to classical Poincaré series, as follows. For any Riemann surface  $R$ , let  $Q(R)$  denote the Banach space of integrable holomorphic

quadratic differentials  $\phi(z)dz^2$  with  $\|\phi\| := \int_R |\phi| < \infty$ . Starting with  $\phi$  in  $Q(\Delta)$ , we can construct an automorphic form for  $G$  by the Poincaré series [Poin]:

$$\Theta(\phi) = \sum_G g^*(\phi).$$

Since  $\Theta(\phi)$  is  $G$ -invariant, it determines an element of  $Q(X)$ .

In Teichmüller theory,  $Q(X)$  is naturally identified with the cotangent space to  $\text{Teich}(X)$  at  $X$ , its norm is dual to the Teichmüller metric, and the operator  $\Theta : Q(\Delta) \rightarrow Q(X)$  is the coderivative  $d\theta^*$ . This gives:

**Corollary 5** (Kra's Theta Conjecture).  $\|\Theta\| < 1$  for classical Poincaré series.

On a global level the theorem says that lifts of Teichmüller mappings can be relaxed (isotoped to mappings of less dilatation):

**Corollary 6.** Let  $f : X_0 \rightarrow X_1$  be a Teichmüller mapping between distinct points in  $\text{Teich}(X)$ . Then the map  $\tilde{f} : \Delta \rightarrow \Delta$  obtain by lifting  $f$  to the universal covers of domain and range is not extremal among quasiconformal with the same boundary values on  $S^1$ .

More generally, these contraction principles apply to a covering  $Y \rightarrow X$  iff the covering is *nonamenable*; see [Mc1].

Now for a typical (acylindrical)  $M^3$  the skinning map  $\sigma$  can be described as follows. Given a Riemann surface  $X$  in  $\text{Teich}(\partial M)$ ,

- (1) form countably many copies of its universal cover  $\tilde{X}$ , then
- (2) glue them together in a pattern determined by the combinatorics of  $M^3$  to obtain a new Riemann surface  $\sigma(X)$ .

The surface  $\sigma(X)$  contains a dense full measure set of open disks each of which is canonically identified with the universal cover of  $X$ . By the results above, step (1) is a contraction for the Teichmüller metric. Step (2) is at worst an isometry, so  $\|d\sigma\| \leq \|d\theta\| < c[X] < 1$ . (A more detailed expository account appears in [Mc5].)

### 3. Renormalization and 3-Manifolds Which Fiber over the Circle

For the special case of Haken manifolds presented as surface bundles over the circle, the construction of a hyperbolic structure is different in spirit and finds parallels with the construction of fixed points for renormalization.

#### 3.1 Surface Bundles

Let  $S$  be a closed oriented surface of genus  $g > 1$ , and let  $\phi : S \rightarrow S$  be a *mapping class*, that is a diffeomorphism determined up to isotopy. From this data one can construct a 3-manifold by starting with  $M^3 = S \times [0, 1]$  and gluing the ends together by  $\phi$ . Every 3-manifold which fibers over the circle admits such a description;  $\phi$  is the *monodromy* of the fibration.

**Theorem 7** [Th2]. *A 3-manifold  $M^3$  which fibers over the circle admits a hyperbolic structure iff the monodromy  $\phi$  is pseudo-Anosov.*

As in the preceding section, the construction of a hyperbolic structure can be formulated as a fixed point problem. There are two essential differences: (1) the desired fixed point lies on the boundary of Teichmüller space, rather than in its interior, and (2) it is dynamically hyperbolic rather than attracting.

**Construction of the Hyperbolic Structure.** We will work in the representation variety  $\mathcal{V} = \text{Hom}(\pi_1(S), \text{Isom}(\mathbb{H}^3))/\text{conjugation}$ . Let  $AH(S) \subset \mathcal{V}$  denote the closed subset of discrete faithful representations. The idea is to find in  $AH(S)$  the  $\mathbb{Z}$ -covering space of  $M^3$  carrying the fundamental group of a fiber. The deck transformation acts by isometry on this covering space, so it is characterized as a fixed point in  $AH(S)$  for the map

$$\Phi : \mathcal{V} \rightarrow \mathcal{V}$$

given by  $\Phi(\varrho) = \varrho \circ \pi_1(\phi)$ .

The construction of the fixed point can be organized into two steps. Let  $QF(S) \subset AH(S)$  be the open subset of quasifuchsian groups; it is holomorphically parameterized by  $\text{Teich}(S) \times \text{Teich}(\bar{S})$  (here  $\bar{S}$  indicates reversal of orientation) and we denote by  $\varrho(X, Y)$  the marked group corresponding to a pair of Riemann surfaces  $X$  and  $Y$ .

**Step 1:** Form the limit  $\varrho_\infty = \lim_{n \rightarrow \infty} \varrho(X, \phi^{-n}(Y))$ . Here  $X$  and  $Y$  are arbitrary Riemann surfaces and  $\phi(Y)$  denotes the action of the mapping class on Teichmüller space. The representations  $\varrho(X, \phi^{-n}(Y))$  range in a Bers' slice, which has compact closure in  $AH(S)$ , so the existence of some accumulation point is clear. Logically one can work with any accumulation point  $\varrho_\infty$ ; in fact, the sequence converges [CT, §7].

**Step 2:** Form the limit  $\varrho = \lim \Phi^n(\varrho_\infty)$ ; this is a fixed point for  $\Phi$ . Existence of this limit depends on

**Compactness:** The double limit theorem of [Th2], which assures there is some accumulation point  $\varrho$ ; and

**Rigidity:** Sullivan's quasiconformal rigidity theorem [Sul1], which gives  $\Phi(\varrho) = \varrho$  for any accumulation point.

The first step produces a point on the stable manifold of a fixed point of  $\Phi$ , and the second iterates it to find the fixed point.

The limit in (Step 2) can be lifted to the level of marked groups  $G_n \subset \text{Aut}(\mathbb{H}^3)$  (rather than groups up to conjugacy) such that  $G_n$  tends algebraically to  $G = \text{Image}(\varrho)$ . The groups  $G_n$  (conjugate to  $\text{Image}(\Phi^n(\varrho_\infty))$ ) are all isomorphic; we are viewing a single dynamical system from a changing perspective. As  $n \rightarrow \infty$  the limit set of  $G_n$  becomes denser and denser, and the limit set of  $G$  is the full sphere.

The group  $G_{n+1}$  is obtained from  $G_n$  by a  $K$ -quasiconformal deformation with uniform  $K$ . By compactness of  $K$ -quasiconformal maps, one obtains a quasiconformal map  $\psi$  equivariant with respect to  $G$  and inducing the automorphism  $\Phi$ . Since the limit set is the whole sphere,  $\psi$  is conformal by Sullivan's result. The group generated by  $G$  and  $\psi$  together is then a Kleinian group isomorphic to  $\pi_1(M^3)$ .

### 3.2 Quadratic-Like Maps

This discussion parallels the emerging complex viewpoint on *renormalization* of quadratic-like maps. For concreteness we will discuss the case of period doubling; see [Cvi, Milnor, Sul3] and [Sul4] for background and more details.

Consider the family of quadratic maps  $z \mapsto z^2 + c$  as the parameter  $c$  decreases along the real axis, starting at  $c = 0$ . One finds a sequence of parameter values  $c(n)$  for which the attractor of  $f_n(z) = z^2 + c(n)$  bifurcates from a cycle of order  $2^n$  to  $2^{n+1}$ ; at  $c(\infty) = \lim c(n)$  the attractor becomes a Cantor set. This *cascade of period doublings* was observed by Feigenbaum to have many universal features around  $f_\infty(z) = z^2 + c(\infty)$ . For example

$$\frac{c(n) - c(\infty)}{c(n+1) - c(\infty)} \rightarrow \lambda = 4.669201609\dots$$

and this value of  $\lambda$  (as well as the fine structure of  $f_\infty$ , such as the Hausdorff dimension of its attracting Cantor set) is the same for other families of smooth mappings with the same topological form as  $z^2 + c$ .

This universality is part of a larger renormalization picture proposed by Feigenbaum and established rigorously by Lanford and others. We present a version with the complex quadratic-like maps of Douady and Hubbard [DH2]; cf. [Sul3].

A *quadratic-like mapping*  $f : U \rightarrow V$  is a proper degree two holomorphic map between open disks with  $\overline{U} \subset V \subset \mathbb{C}$ . Its *filled-in Julia set*  $K(f)$  is  $\bigcap_1^\infty f^{-n}(V)$ . When  $K(f)$  is connected, there is a unique quadratic polynomial  $I(f)$  (the *inner class*) conjugate to  $f$  near  $K(f)$  by a quasiconformal map which is conformal on  $K(f)$ . Thus  $I$  takes values in the Mandelbrot set  $M$  of polynomials  $z^2 + c$  with connected Julia sets.

Let  $\mathcal{Q}$  be the space of all analytic maps  $f : \Omega_f \rightarrow \mathbb{C}$  defined on a region  $\Omega_f$  containing the origin, such that  $f'(0) = 0$  and  $f$  is quadratic-like on some neighborhood of zero. We identify  $f(z)$  and  $g(z)$  if some rescaling  $\alpha f(z/\alpha)$  agrees with  $g(z)$  on their common domain of definition. Finally  $f_i \rightarrow f$  if there are representatives of  $f_i$  which converge to  $f$  uniformly on compact subsets of  $\Omega_f$ .

The *renormalization operator*  $\mathcal{R} : \mathcal{D}' \rightarrow \mathcal{Q}$  is given by  $\mathcal{R}(f) = f \circ f$ . It is defined on an open set  $\mathcal{D}'$  such that  $f \circ f$  is still quadratic-like near the origin.

Central to the picture is the existence of a unique fixed point  $g$  for  $\mathcal{R}$ . Since  $g$  and  $\mathcal{R}(g)$  are equivalent,  $g$  satisfies the Cvitanović-Feigenbaum functional equation  $\alpha g \circ g(z) = g(\alpha z)$ . The universal constant  $\lambda$  above is the unique expanding eigenvalue for  $\mathcal{R}$  at  $g$ .

We will sketch a construction of this fixed point  $g$  which parallels the geometrization of surface bundles.

Douady and Hubbard define a *tuning map*  $\tau : M \rightarrow M$  such that  $I(\mathcal{R}(f)) = \tau^{-1}(I(f))$  when defined. Thus  $\tau$  describes the inverse of renormalization as it acts on the inner class. One finds that  $\tau(f_n) = f_{n+1}$  and  $\tau$  fixes the Feigenbaum polynomial  $f_\infty$ . This accomplishes:

*Step 1:* Form the limit  $f_\infty = \lim \tau^n(f_0)$ .

Now let  $\mathcal{Q}_\infty$  denote those  $f$  with inner class  $I(f) = f_\infty$ . Then  $\mathcal{R}(\mathcal{Q}_\infty) \subset \mathcal{Q}_\infty$  and  $f$  and  $\mathcal{R}(f)$  are quasiconformally conjugate near  $K(f)$  for any  $f$  in  $\mathcal{Q}_\infty$ .

*Step 2:* Form the limit  $g = \lim \mathcal{R}^n(f_\infty)$ . This is a fixed point for  $\mathcal{R}$ , and in fact all  $f$  in  $\mathcal{Q}_\infty$  are attracted to  $g$ .

The proof of (Step 2) again appeals to two principles.

**Compactness:** For any  $f$  in  $\mathcal{Q}_\infty$ ,  $\langle \mathcal{R}^n(f) \rangle$  ranges in a compact subset of  $\mathcal{Q}$ . This is a fundamental result of Sullivan [Sul4]. Thus there is a subsequence of  $n$  such that  $\mathcal{R}^n(f) \rightarrow g_0$ ,  $\mathcal{R}^{n-1}(f) \rightarrow g_1$ , ...,  $\mathcal{R}^{n-k}(f) \rightarrow g_k$  and the tower  $\langle g_0, g_1, \dots \rangle$  satisfies  $\mathcal{R}(g_k) = g_{k-1}$ . We can then apply:

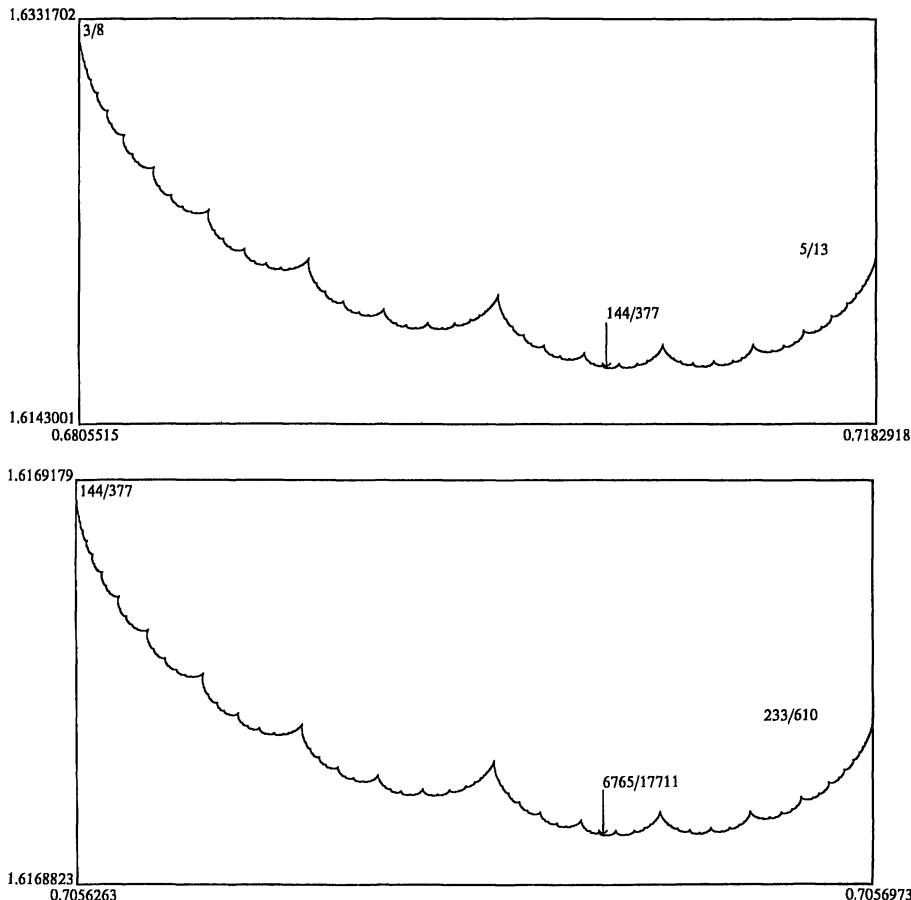
**Rigidity:** Such a tower admits no quasiconformal deformations [Mc4].

Since  $\langle g_0, g_1, \dots \rangle$  is conjugate to  $\langle g_1, g_2, \dots \rangle$  by a suitable limit of a quasiconformal conjugacy between  $f$  and  $\mathcal{R}(f)$ , these towers are conformally identical, and in particular  $\mathcal{R}(g_0) = g_0$ . By rigidity of all limiting towers, the full sequence  $\mathcal{R}^n(f) \rightarrow g_0$  and this fixed point is unique.

**Geometric Limits.** The fixed point of renormalization  $g$  is *not* itself rigid. Its universal structure is a result of being embedded deep in the dynamics of  $f$ . The tower  $\langle g_0, g_1, \dots \rangle$  can be thought of as a geometric limit of the dynamical system generated by  $f$  as one rescales about its critical point. The limiting dynamic is *divisible* ( $g_n = g_{n+1} \circ g_{n+1}$ ), and its Julia set fills the whole plane. If we set  $g_{-n}$  = the  $2^n$ th iterate of  $g_0$ , then renormalization acts as a shift on the bi-infinite tower  $\langle \dots, g_{-1}, g_0, g_1, \dots \rangle$  in a manner reminiscent of the deck transformation acting on the  $\mathbb{Z}$ -covering space of  $M^3$  constructed before.

**Self-Similarity in the Mandelbrot Set and in Bers' Slice.** In the polynomial-like setting, one can actually define countably many renormalization and tuning operators  $\mathcal{R}_c$ ,  $\tau_c$ , one for each  $c$  such that the critical point of  $z^2 + c$  is periodic. Milnor has made a detailed computer study of these operators, supporting many conjectures [Milnor]; among them, that  $\tau_c$  has a unique fixed point and is differentiable there, with derivative given by the inverse of the expanding eigenvalue of  $\mathcal{R}_c$  at its fixed point.

Similarly, we conjecture (in the case of one dimensional Teichmüller spaces) that Bers' boundary is self-similar about the point  $g_\infty$  constructed in (Step 1), with similarity factor given by the expanding eigenvalue of the mapping class  $\Phi$  at the fixed point of (Step 2). Dave Wright's computer study of the closely related Maskit boundary for the Teichmüller space of a punctured torus supports this conjecture [Wr]. In this case the mapping class group is  $\text{SL}_2\mathbb{Z}$ , and the expanding



**Fig. 1.** Self-similarity at the edge of Teichmüller space

eigenvalue is algebraic (but different from the eigenvalues of the matrix). For example, when  $\Phi = \begin{pmatrix} 21 \\ 11 \end{pmatrix}$ , the boundary scales by  $\lambda = 4.79129\dots = \frac{5+\sqrt{21}}{2}$ ; see Fig. 1 for two blowups around  $\varrho_\infty$  computed by Wright.

**Remark.** Sullivan has established a compactness theorem for arbitrary compositions of a finite number of renormalization operators  $\mathcal{R}_c$ , with the condition that  $c$  is real. Our rigidity argument applies whenever such a compactness result is available. Much progress on a conceptual understanding of the full renormalization picture, including a different approach to rigidity, appears in [Sul4].

#### 4. Boundaries and Laminations

We conclude with a very brief account of progress on the boundary of Teichmüller space and the boundary of the Mandelbrot set.

**Conjecture** (Douady–Hubbard). The Mandelbrot set  $M$  is locally connected. Its boundary is homeomorphic to a quotient of the circle by an explicit combinatorial equivalence relation. (Cf. [DH1, Dou, Lav, Th3])

**Conjecture** (Thurston). Bers' boundary for Teichmüller space, modulo quasiconformal equivalence, is homeomorphic to the space  $\mathbb{PML}$  of projective measured laminations, modulo forgetting the measure.

Both conjectures express the hope that certain geometrically infinite dynamical systems can be uniquely described by a lamination on the circle — invariant under  $z^2$  in the first case, and under the action of a surface group in the second.

Measures supported on maximal systems of disjoint simple closed curves are dense in  $\mathbb{PML}$ ; these correspond to *maximal cusps* in Bers' boundary, that is geometrically finite limits of quasifuchsian groups where these curves have been pinched to form rank one cusps. Thus Thurston's conjecture is supported by:

**Theorem 8** [Mc3]. *Maximal cusps are dense in Bers' boundary.*

This result was conjectured by [Bers]. The proof is by via an explicit estimate for the algebraic effect of a quasiconformal deformation supported in the thin part.

Also relevant is Bonahon's result: a general geometrically infinite surface group admits an ending lamination [Bon], supporting the conjecture that geometrically finite groups are dense in  $AH(S)$ .

Progress on the Mandelbrot set includes the following breakthrough:

**Theorem 9** (Yoccoz).  *$M$  is locally connected at every quadratic polynomial which is not in the image of a tuning map.*

Yoccoz's result brings us a step closer to resolving the well-known:

**Conjecture.** Hyperbolic dynamics is open and dense in the space of complex quadratic polynomials.

It seems likely that Yoccoz's theorem generalizes to the case of polynomials lying in the image of only finitely many tuning maps. If so, by [MSS], the density of hyperbolic dynamics is equivalent to the quasiconformal rigidity of infinitely renormalizable polynomials.

## References

- [Bers] L. Bers: On boundaries of Teichmüller spaces and on kleinian groups; I. Ann. Math. **91** (1970) 570–600
- [Bon] F. Bonahon: Bouts des variétés hyperboliques de dimension 3. Ann. Math. **124** (1986) 71–158
- [CT] J. W. Cannon, W. P. Thurston: Group invariant Peano curves. Preprint.
- [Cvi] P. Cvitanović: Universality in chaos. Adam Hilger Ltd, 1984.
- [Dou] A. Douady: Algorithms for computing angles in the Mandelbrot set. In: M. F. Barnsley, S. G. Demko, (eds.) Chaotic dynamics and fractals. Academic Press, 1986, pp. 155–168

- [DH1] A. Douady, J. Hubbard: Étude dynamique des polynômes complexes. Pub. Math. d'Orsay, 1984
- [DH2] A. Douady, J. Hubbard: On the dynamics of polynomial-like mappings. Ann. Sci. Éc. Norm. Sup. **18** (1985) 287–344
- [DH3] A. Douady, J. Hubbard: A proof of Thurston's topological characterization of rational maps. To appear in Acta Mathematica
- [Lav] P. Lavaurs: Une description combinatoire de l'involution définie par  $M$  sur les rationnels à dénominateur impair. CRAS Paris **303** (1986) 143–146
- [Mc1] C. McMullen: Amenability, Poincaré series and quasiconformal maps. Invent. math. **97** (1989) 95–127
- [Mc2] C. McMullen: Iteration on Teichmüller space. Invent. math. **99** (1990) 425–454
- [Mc3] C. McMullen: Cusps are dense. Ann. Math. **133** (1991) 217–247
- [Mc4] C. McMullen: Renormalization and 3-manifolds which fiber over the circle. In preparation
- [Mc5] C. McMullen: Riemann surfaces and the geometrization of 3-manifolds. Submitted to the Annals of the New York Academy of Sciences
- [Milnor] J. Milnor: Self-similarity and hairiness in the Mandelbrot set. In M. C. Tangora (ed.) Computers in geometry and topology. Lect. Notes Pure Appl. Math. Dekker, 1989
- [Mor] J. Morgan: On Thurston's uniformization theorem for three dimensional manifolds. In: The Smith conjecture. Academic Press, 1984, pp. 37–125
- [Mum] D. Mumford: A remark on Mahler's compactness theorem. Proc. AMS **28** (1971) 289–294
- [MSS] R. Mañé, P. Sad, D. Sullivan: On the dynamics of rational maps. Ann. Sci. Éc. Norm. Sup. **16** (1983) 193–217
- [Poin] H. Poincaré: Mémoire sur les fonctions Fuchsiennes. Acta Math. **1** (1882/3), 193–294
- [Sul1] D. Sullivan: On the ergodic theory at infinity of an arbitrary discrete group of hyperbolic motions. In: Riemann surfaces and related topics: Proceedings of the 1978 Stony Brook Conference. Ann. Math. Studies 97. Princeton, 1981
- [Sul2] D. Sullivan: Quasiconformal homeomorphisms and dynamics I: Solution of the Fatou-Julia problem on wandering domains. Ann. Math. **122** (1985) 401–418
- [Sul3] D. Sullivan: Quasiconformal homeomorphisms in dynamics, topology and geometry. In: Proceedings of the International Conference of Mathematicians. American Math. Soc., 1986, pp. 1216–1228
- [Sul4] D. Sullivan: Bounds, quadratic differentials and renormalization conjectures. In: American Mathematical Society Centennial Publications, volume 2: Mathematics into the twenty-first century. Amer. Math. Soc., to appear
- [Th1] W. P. Thurston: Hyperbolic structures on 3-manifolds I: Deformations of acylindrical manifolds. Ann. Math. **124** (1986) 203–246
- [Th2] W. P. Thurston: Hyperbolic structures on 3-manifolds II: Surface groups and 3-manifolds which fiber over the circle. To appear in Ann. Math.
- [Th3] W. P. Thurston: On the combinatorics and dynamics of iterated rational maps. Preprint.
- [Wr] D. Wright: The shape of the boundary of Maskit's embedding of the Teichmüller space of once-punctured tori. Preprint.



# Analytic Capacity for Arcs

Takafumi Murai

Department of Mathematics, School of Science, Nagoya University, Nagoya, 464-01 Japan

## 1. Introduction

For a domain  $\Omega$  in the extended complex plane  $\mathbb{C} \cup \{\infty\}$ ,  $H^\infty(\Omega)$  denotes the Banach space of bounded analytic functions in  $\Omega$  with supremum norm  $\|\cdot\|_{H^\infty}$ . For  $\zeta \in \Omega$ , we put  $c(\zeta; \Omega) = \sup |f'(\zeta)|$ , where the supremum is taken over all  $f \in H^\infty(\Omega)$ ,  $\|f\|_{H^\infty} \leq 1$ , and  $f'(\infty)$  is defined by  $\lim_{z \rightarrow \infty} z\{f(\infty) - f(z)\}$ . Given  $f \in H^\infty(\Omega)$ ,  $\|f\|_{H^\infty} \leq 1$ , we take  $g(z) = \{f(z) - f(\zeta)\}/\{1 - f(z)\overline{f(\zeta)}\}$ . Then  $\|g\|_{H^\infty} \leq 1$ ,  $g(\zeta) = 0$  and  $|g'(\zeta)| \geq |f'(\zeta)|$ . Thus, for the computation of  $c(\zeta; \Omega)$ , we may restrict our attention to functions vanishing at  $\zeta$ . For  $p \geq 1$  and a domain  $\Omega$  bounded by a finite number of analytic Jordan curves,  $H^p(\Omega)$  denotes the  $H^p$ -space of analytic functions  $f$  in  $\Omega$  with norm  $\|f\|_{H^p} = \{(1/2\pi) \int_{\partial\Omega} |f'|^p |dz|\}^{1/p}$ , where  $\partial\Omega$  is the boundary of  $\Omega$ . The condition  $\int_{\partial\Omega} (\mathcal{M}f)^p |dz| < \infty$  is required for each  $f \in H^p(\Omega)$ , where  $\mathcal{M}f$  is a non-tangential maximal function which controls the behaviour of  $f$  near the boundary. The analytic capacity of a compact set  $E$  in  $\mathbb{C}$  is defined by  $\gamma(E) = c(\infty; \Omega_E)$ , where  $\Omega_E$  is the component of  $E^c$  containing  $\infty$ . Analytic capacity plays an important role in the theory of conformal mapping [SO], the 2-dimensional fluid dynamics [Mi], approximation theory [Ga, V2, Z] and singular integrals [Ch, J, Mu3]. Ahlfors [A] shows that a compact set  $E$  satisfying Riemann's theorem on removable singularities is characterized by  $\gamma(E) = 0$ , and Garabedian [G] studies  $\gamma(\cdot)$  from the point of view of the dual extremum problem. Basic properties are mentioned in [AB, Ga, Z]. Here is the classical method of the computation of  $\gamma(E)$  in the case where  $E^c = \Omega_E$  and  $\partial E$  consists of a finite number of analytic Jordan curves.

The Ahlfors-Garabedian method [A, G]: *Construct a pair  $(f_0, \psi_0)$  of functions so that  $f_0 \in H^\infty(E^c)$ ,  $\|f_0\|_{H^\infty} \leq 1$ ,  $f_0(\infty) = 0$ ,  $\psi_0 \in H^1(E^c)$ ,  $\psi_0(\infty) = 1$  and*

$$(1) \quad \frac{1}{i} \int_{\partial E} f_0 \psi_0 \, dz = |\psi_0| |dz| \text{ almost everywhere (a.e.) on } \partial E,$$

where the orientation of  $dz$  is chosen so that  $E^c$  lies to the left. Once such a pair has been constructed, we have

$$(2) \quad \gamma(E) = f'_0(\infty) = \|\psi_0\|_{H^1}.$$

In fact,

$$\begin{aligned}
|f'_0(\infty)| &\leq \gamma(E) = \sup_0 |f'(\infty)| = \sup_0 \left| \frac{1}{2\pi} \int_{\partial E} f \, dz \right| \\
&= \sup_0 \left| \frac{1}{2\pi} \int_{\partial E} f \psi_0 \, dz \right| \quad (\text{by } \psi_0(\infty) = 1) \\
&\leq \|\psi_0\|_{H^1} = \frac{1}{2\pi i} \int_{\partial E} f_0 \psi_0 \, dz \quad (\text{by (1)}) \\
&= \frac{1}{2\pi i} \int_{\partial E} f_0 \, dz = f'_0(\infty),
\end{aligned}$$

which yields (2). Here  $\sup_0$  is the supremum over all  $f \in H^\infty(E^c)$ ,  $\|f\|_{H^\infty} \leq 1$ ,  $f(\infty) = 0$ . Thus it is essential to construct the pair  $(f_0, \psi_0)$ . For any compact set  $F$ , there exists uniquely  $f(\cdot; \Omega_F) \in H^\infty(\Omega_F)$  with norm 1 such that  $f'(\infty; \Omega_F) = \gamma(F)$  [Ga, p. 24]. This is called the Ahlfors function of  $F$ . Let  $\{\Omega_n\}_{n=1}^\infty$  be an increasing sequence of domains  $\ni \infty$  with smooth boundaries such that  $\bigcup_{n=1}^\infty \Omega_n = \Omega_F$ . Then there exists uniquely  $\psi_n \in H^1(\Omega_n)$  such that  $\psi_n(\infty) = 1$  and  $\|\psi_n\|_{H^1} = \gamma(\Omega_n^c)$ . The sequence  $\{\psi_n\}_{n=1}^\infty$  converges in  $\Omega_F$  and the limit  $\psi(\cdot; \Omega_F)$  is called the Garabedian function of  $F$ ;  $\psi(\cdot; \Omega_F)$  is determined independent of  $\{\Omega_n\}$  [Sm, Su1, 2]. Note that  $f_0 = f(\cdot; E^c)$  and  $\psi_0 = \psi(\cdot; E^c)$ . Our main theme is the study of  $\gamma(\cdot)$ ,  $f(\cdot; \cdot)$  and  $\psi(\cdot; \cdot)$ . The study from the point of view of Green's functions and harmonic measures is classical. In this talk, we focus on an approach based on the singular integral operator  $\mathcal{H}$  defined by Calderón [C]. Let  $\mathcal{A}$  denote the totality of sets  $E$  such that  $E$  consists of a finite number of mutually disjoint closed arcs  $\{C_j\}_{j=1}^n$  and each arc  $C_j$  is a finite union of analytic arcs. For the computation of  $\gamma(\cdot)$ , we restrict our attention to sets in  $\mathcal{A}$ . Here are two motivations to deal with sets in  $\mathcal{A}$ . For a compact set  $F$  with a smooth boundary, we can discuss the Hadamard variation and the Schiffer variation [S2] of  $\gamma(F)$ . Then we can express  $\gamma(F)$  as a perturbation from  $\gamma(E)$  for a set  $E$  in  $\mathcal{A}$ . Thus, *in order to get global properties of  $\gamma(\cdot)$* , it is necessary to study  $\gamma(E)$ ,  $E \in \mathcal{A}$ . Another motivation is as follows. Given a compact set  $F$ , we can find a finite union  $G$  of closed disks such that  $|\gamma(F) - \gamma(G)|$  is arbitrarily small. Note that  $\gamma(G) = \gamma(\partial G)$ . Removing some arcs on  $\partial G$ , we can find  $E \in \mathcal{A}$  such that  $|\gamma(G) - \gamma(E)|$  is arbitrarily small. (In this approximation,  $E$  can be chosen so that the connectivity of  $E$  is less than or equal to that of  $F$ .) Thus it is sufficient to study  $\gamma(E)$ ,  $E \in \mathcal{A}$ . For  $E \in \mathcal{A}$ ,  $\partial E$  denotes the boundary of  $E$  having two sides. The  $H^p$ -space  $H^p(E^c)$  of analytic functions in  $E^c$  is analogously defined as above.

## 2. The Singular Integral Operator $\mathcal{H}$

For  $E \in \mathcal{A}$ ,  $L^2(E)$  denotes the  $L^2$ -space of functions on  $E$  with respect to the arc-length  $|dz|$ . The singular integral operator  $\mathcal{H}_E$  from  $L^2(E)$  to itself is defined by

$$\mathcal{H}_E h(z) = \frac{1}{\pi} \text{p.v.} \int_E \frac{1}{\zeta - z} h(\zeta) |d\zeta|,$$

where p.v. is the principal value. There are many articles about  $\mathcal{H}$  [C, Ch, CJS, CMM, D, Me, Mu3]. The operator  $\bar{\mathcal{H}}_E$  is defined by  $\bar{\mathcal{H}}_E h = \mathcal{H}_E \bar{h}$  and the inverse operator of  $\text{Id} - \mathcal{H}_E \bar{\mathcal{H}}_E$  is denoted by  $\mathcal{T}_E$ , where  $\text{Id}$  is the identity operator. Here are fundamental expressions of  $\psi(\cdot; E^c)$ ,  $f(\cdot; E^c)$  and  $\gamma(E)$  in terms of  $\mathcal{H}_E$  [Mu4]:

$$\begin{aligned}\psi(z; E^c) &= \left\{ 1 + \frac{1}{\pi} \int_E \frac{1}{\zeta - z} \bar{\mathcal{H}}_E \mathcal{T}_E 1(\zeta) |d\zeta| \right\}^2, \\ f(z; E^c) &= \frac{1}{\pi} \int_E \frac{1}{\zeta - z} \bar{\mathcal{T}}_E 1(\zeta) |d\zeta| / \sqrt{\psi(z; E^c)}, \\ (3) \quad \gamma(E) &= \frac{1}{\pi} \int_E \mathcal{T}_E 1(\zeta) |d\zeta|.\end{aligned}$$

The proof is founded on Garabedian's duality theorem [G]:  $\gamma(E) = \inf\{\|\phi\|_{H^2}^2; \phi \in H^2(E^c), \phi(\infty) = 1\}$ . This shows that

$$\gamma(E) = \inf \left\{ \frac{1}{\pi} \int_E (|1 + \mathcal{H}_E h|^2 + |h|^2) |dz|; h \in L^2(E) \right\}.$$

Using the standard variational method, we obtain the required formulae. The following formula plays an important role to compute  $\gamma(E)$  practically:

$$(4) \quad \mathcal{H}_E \{u \mathcal{H}_E v + \mathcal{H}_E u \cdot v\} = \mathcal{H}_E u \cdot \mathcal{H}_E v - \varrho_E u v \quad (u, v \in L^2(E), \varrho_E = \overline{dz/dz}).$$

The other relations between  $\gamma(E)$  and  $\mathcal{H}_E$  are mentioned in [Mu3]; roughly speaking,  $\gamma(E)$  is comparable with  $1/\|\mathcal{H}_E\|_{1,w}$ , where  $\|\mathcal{H}_E\|_{1,w}$  denotes the norm of  $\mathcal{H}_E$  as an operator from the  $L^1$  space to the weak  $L^1$  space. If  $E$  is contained in the real line  $\mathbb{R}$ ,  $\mathcal{H}_E$  is called the Hilbert transform and denoted by  $H_E$ . Formula (4) yields that

$$\begin{aligned}\mathcal{T}_E h &= \frac{1}{2} h + \frac{1}{4} \{\tau_E H_E(h\tau_E^{-1}) - \tau_E^{-1} H_E(h\tau_E)\}, \\ H_E \mathcal{T}_E h &= \frac{1}{4} \{\tau_E H_E(h\tau_E^{-1}) + \tau_E^{-1} H_E(h\tau_E)\} \quad (h \in L^2(E), E \subset \mathbb{R}),\end{aligned}$$

where  $\tau_E = \exp\{(\pi/4)H_E 1\}$ . Thus  $\mathcal{T}_E$  and  $H_E \mathcal{T}_E$  handle easily in the case of  $E \subset \mathbb{R}$ . As application of our method, we obtain the following formulae:

$$\begin{aligned}\gamma(E) &= |E|/4 \quad (E \subset \mathbb{R}) \quad [\text{P}], \quad \gamma(E) = \sin(|E|/4) \quad (E \subset \mathbb{T}), \\ \delta(E; 0, \infty) &= 2 \tan(|E|/8) \quad (E \subset \mathbb{T}) \quad [\text{Mu6}].\end{aligned}$$

Here  $\mathbb{T}$  is the unit circle,  $|\cdot|$  is the 1-dimensional Lebesgue measure and  $\delta(E; 0, \infty)$  denotes the supremum of  $|f(0) - f(\infty)|$  over all  $f \in H^\infty(E^c)$ ,  $\|f\|_{H^\infty} \leq 1$ . Using (4), we can find the concrete forms of  $(f(\cdot; \cdot), \psi(\cdot; \cdot))$  in various cases; once a pair has been found, the check is very easy as stated in the introduction.

### 3. Null Sets

The 1-dimensional Hausdorff measure is also denoted by  $|\cdot|$  [F, p. 7]. We have  $\gamma(E) \leq |E|/\pi$  (Painlevé). This shows that  $\gamma(E) = 0$  if the Hausdorff dimension

[F, p. 7]  $\dim(E)$  is less than 1. On the other hand,  $\gamma(E) > 0$  if  $\dim(E) > 1$ . Thus the case of  $\dim(E) = 1$  is critical. Vitushkin [V1] constructed a planar Cantor set  $P_\infty$  such that  $\gamma(P_\infty) = 0$  and  $|P_\infty| > 0$  (cf. [Ga, p. 87]). We are interested in the geometric structure of sets of Vitushkin-Garnett type. Here is a deformation of  $P_\infty$  which handles easily:  $Q_\infty = \bigcap_{n=0}^{\infty} \{\bigcup_{k=k}^{\infty} Q_k\}^\circ$ , where  $Q_0 = [0, 1]$  and

$$Q_k = \left\{ x + i \sum_{j=1}^k (-1)^j 2^{-j} \operatorname{sign}(\sin \pi 2^{j-1} x); 0 \leq x \leq 1 \right\}^\circ \quad (k \geq 1).$$

(The notation  $E^\circ$  denotes the closure of  $E$  and  $\operatorname{sign} 0 = 1$ .) Then  $\gamma(Q_\infty) = 0$  and  $|Q_\infty| > 0$ . There are two methods of the proof of  $\gamma(Q_\infty) = 0$ . The first method is as follows: *Supposing that the nontrivial Ahlfors function of  $Q_\infty$  exists, show a contradiction* [J, M1]. Mattila applies Besicovitch's set theory [F], and Jones uses the BMO norm  $\{(i/2) \sup_{z \in \Omega} \int \int_Q G(z, w) |f'(w)|^2 dw \wedge d\bar{w}\}^{1/2}$  ( $Q = Q_\infty^c$ ), where  $G(z, w)$  is the Green's function of  $\Omega$ . The second method is based on the construction of the approximate Garabedian function [Mu1]: *For  $\varepsilon > 0$ , construct a pair  $(\psi, R)$  of a function  $\psi \in H^1(R)$  and a domain  $R$  so that  $\infty \in R \subset Q_\infty^c$ ,  $\psi(\infty) = 1$  and  $\|\psi\|_{H^1} \leq \varepsilon$ .* Once such a pair has been constructed, we have  $\gamma(Q_\infty) \leq \gamma(R^c) \leq \|\psi\|_{H^1} \leq \varepsilon$ , which gives  $\gamma(Q_\infty) = 0$ . The following dipole function plays an important role in the construction of such a pair:

$$p(z) = \exp \left\{ \frac{\xi}{z-b} - \frac{\xi}{z-a} \right\} \quad \left( a, b \in \mathbb{C}, \frac{\xi}{b-a} > 0 \right).$$

We have  $p(\infty) = 1$ ,  $p(z) = 1 + O(|z|^{-2})$  ( $z \rightarrow \infty$ ) and  $|p(z)| < 1$  in the strip with width  $|b-a|$  which is perpendicular to the segment with endpoints  $a, b$  and contains it. The required function  $\psi$  is expressed as a product of dipole functions. This method works for non-homogeneous Cantor sets. Next we show some estimates of  $\gamma(\cdot)$  from below Calderón, Havin and Marshall [Ma] show that  $\gamma(E) > 0$  if  $E$  is a compact set on a rectifiable curve satisfying  $|E| > 0$ ; this theorem was formerly called the Denjoy conjecture [Ga, p. 36] and, in the proof, Calderón's theorem [C] on  $\mathcal{H}$  plays an important role. Let  $Bu(\cdot)$  denote the Buffon needle probability (the Favard length). Since  $\gamma(\cdot) \leq |\cdot|/\pi$  and  $Bu(\cdot) \leq \text{Const}|\cdot|$ , it is interesting to compare  $\gamma(\cdot)$  with  $Bu(\cdot)$  (Vitushkin's problem [Ma, V2]). Jones and Murai [JM] show that there exists a compact set  $E$  such that  $Bu(E) = 0$  and  $\gamma(E) > 0$  (cf. [M2, Mu2, 3]). This theorem suggests that  $Bu(\cdot) \leq \text{Const} \gamma(\cdot)$ . There are many problems about this topic [HHN, pp. 485–514]. We here note two problems: (I) *Suppose that  $|E| < \infty$ . Does  $Bu(E) = 0$  imply  $\gamma(E) = 0$*  [HHN, p. 491]? (II) *Construct a compact set  $E$  so that  $\gamma_0(E) = 0$  and  $\gamma(E) > 0$*  [Ga, p. 55]. Here  $\gamma_0(E)$  denotes the supremum of  $|\int d\mu|$  over all Cauchy potentials  $\mathcal{C}\mu(z) = \int \frac{1}{\zeta-z} d\mu$  of measures  $\mu$  on  $E$  such that  $\|\mathcal{C}\mu\|_{H^\infty} \leq 1$ . Problem

I originates in Besicovitch's set theory. A Borel set  $E$  satisfying  $|E| < \infty$  is regular if  $d(z, E) = 1$   $|\cdot|$ -a.e. on  $E$ , and  $E$  is irregular if  $d(z, E) < 1$   $|\cdot|$ -a.e. on  $E$ , where  $d(z, E) = \liminf_{r \rightarrow 0} |\{\zeta \in E; |\zeta - z| \leq r\}|/(2r)$ . If  $|E| < \infty$ , then  $Bu(E) = 0$  is equivalent to the irregularity of  $E$  [F, p. 89]. Problem II is posed to clarify the difference between {bounded Cauchy potentials} and  $H^\infty(\cdot)$ . Take a bad arc  $E$  (like a snowflake) such that  $\dim(E) = 1$  and the diameter is equal to 1. Note that

$\gamma(E) \geq 1/4$  [Ga, p. 9]. Is  $\gamma_0(E)$  small? Study Ahlfors functions which cannot be expressed as Cauchy potentials of measures.

## 4. Projection

Let  $\mathcal{L}_\theta$  ( $-\pi/2 < \theta \leq \pi/2$ ) denote the straight line  $x \sin \theta = y \cos \theta$  and let  $\text{pr}_\theta E$  denote the projection of  $E$  to  $\mathcal{L}_\theta$ . As is well known, a regular set is contained in a countable union of rectifiable graphs [F, p. 45]. Thus CHM's theorem [Ma] shows that  $\gamma(E) > 0$  if  $E$  is a regular set satisfying  $|E| > 0$ . A set  $E$  satisfying  $|E| < \infty$  is irregular if  $|\text{pr}_\theta E| = |\text{pr}_{\theta'} E| = 0$  for two distinct numbers  $\theta, \theta'$  [F, p. 90]. From this point of view, it is interesting to estimate  $\gamma(E)$  in terms of the projection of  $E$  to one direction. Recall that  $\gamma(Q_\infty) = 0$  and  $\text{pr} Q_\infty = [0, 1]$ , where  $\text{pr} = \text{pr}_0$ . To understand the geometric meaning of  $Q_\infty$ , we begin with  $Q_1$ . Let  $\Gamma(z) = [-1/2, 1/2] \cup (z + [-1/2, 1/2])$  and  $\gamma(z) = \gamma(\Gamma(z))$ . In hydrodynamics,  $\Gamma(z)$  is regarded as a biplane wing section [Mi, Chap. VII] and there are many articles about  $\Gamma(z)$  [Fe, Gar]. In order to practically compute  $\gamma(z)$ , we introduce the lift coefficient  $\mathcal{L}(z)$  of  $\Gamma(z)$  [Fe, Mi, p. 203]. There exists uniquely  $f_z \in H^1(\Gamma(z)^c)$  such that  $f_z$  is real-valued continuous on  $\partial\Gamma(z) - \{\pm 1/2, z \pm 1/2\}$ ,  $f_z(\infty) = -i$  and  $f_z$  satisfies Joukowski's hypothesis [Mi, p. 199] " $|f_z(\zeta)| < \infty$  ( $\zeta = 1/2, z + 1/2$ )". Taking account of Blasius' formula [Mi, p. 173], we define the lift coefficient  $\mathcal{L}(z)$  of  $\Gamma(z)$  by

$$\mathcal{L}(z) = \frac{1}{4} \left| \frac{1}{2\pi} \int_{\partial\Gamma(z)} f_z(\zeta)^2 d\zeta \right| \left( = \frac{1}{2} |f_z'(\infty)| \right).$$

It is sufficient to study  $\gamma(z)$  in  $\mathbb{P} = \{\text{Re } z \geq 0, \text{Im } z \geq 0\} - [0, 1]$ . Then the following assertion [Mu5] holds:

$$(5) \quad \gamma(z) = \frac{1}{2} + \frac{\text{Im } z}{2} \int_{\lambda(z)} \left\{ \frac{\gamma(\zeta)}{\mathcal{L}(\zeta)} - 1 \right\} \frac{d(\text{Im } \zeta)}{(\text{Im } \zeta)^2} \quad (z \in \mathbb{P}),$$

where  $\lambda(z)$  denotes the arc in  $\mathbb{P}$  with endpoints  $z$  and a positive number such that (the modulus [SO, p. 199] of  $\Gamma(\zeta)^c$ ) is invariant on  $\lambda(z)$  and,  $z$  is chosen as the starting point of the curvilinear integral. The inequality  $\mathcal{L}(z) \leq \gamma(z)$  holds, and the equality  $\mathcal{L}(z) = \gamma(z)$  holds if and only if  $\text{Im } z = 0$ . Using (5), we obtain the following equality [Mu5]:

$$\min_{x \geq 0, y \geq 0} \gamma(x + iy)/\gamma(x) = \min_{y \geq 0} \gamma(1 + iy)/\gamma(1)(0.9\dots).$$

Note that  $\gamma(x) = |\Gamma(x)|/4|\Gamma(x)|/4$  ( $x \in \mathbb{R}$ ) and the projection of  $\Gamma(1 + iy)$  to  $\mathbb{R}$  overlaps only at 1/2. This is a reason why we take  $Q_1$  as the first step to construct a compact set  $Q_\infty$  of Vitushkin-Garnett type. Even  $\Gamma(z)$ , the behaviour of  $\gamma(z)$  is not simple. If  $0 < x_0 < 1$  is sufficiently near to 1, then  $\gamma(x_0 + iy)$  has at least two local extrema in  $(0, \infty)$  as a function of  $y$ . From this fact, we conjecture that there exists a compact set  $E$  such that  $\gamma(E) = 0$  and  $\gamma(T(E)) > 0$ , where  $T(x + iy) = x + i2y$  (or  $= x + (iy/2)$ ). The following theorem gives the geometric information about compact sets  $E$  such that  $\gamma(E)/|\text{pr } E|$  is small.

**Theorem 1** [Mu2-4]. *If  $E$  is a compact set on a graph  $\Gamma$  such that  $|\Gamma| \leq 1$ , then  $\gamma(E) \geq C|\text{pr } E|^{3/2}$ , where  $C$  is an absolute constant. The power  $3/2$  is best possible.*

This is a generalization of CHM's theorem with a quantitative estimate and an interpretation, in terms of  $\gamma(\cdot)$ , of the optimal estimate of  $\|C[\cdot]\|$  defined later. Dilating the coordinate axes, we obtain  $\gamma(E) \geq C|\text{pr } E|^{3/2}|\Gamma|^{-1/2}$  for any compact set  $E$  on a rectifiable graph  $\Gamma$  containing  $E$ . Thus, if  $|\text{pr } E| \geq 1$  and  $|\Gamma| \leq M$ , then  $\gamma(E) \geq C/\sqrt{M}$ . We have  $\gamma(Q_n) \geq \text{Const}/\sqrt{n}$ , for example. The first half assertion in this theorem is rewritten in the following form also: *If  $\gamma(E) \leq \varepsilon$ , then a graph of length less than  $C^2\delta^3\varepsilon^{-2}$  does not contain any subset  $F$  of  $E$  satisfying  $|\text{pr } F| \geq \delta$ .* In order to prove the inequality in Theorem 1, it is necessary to investigate  $\mathcal{H}$  in detail. Let BMO denote the Banach space of functions on  $\mathbb{R}$ , modulo constants, of bounded mean oscillation. For a real-valued function  $a \in \text{BMO}$ , the singular integral operator  $C[a]$  from the  $L^2$ -space of functions on  $\mathbb{R}$  to itself is defined by

$$C[a]h(x) = \frac{1}{\pi} \text{p.v.} \int_{-\infty}^{\infty} \frac{1}{y-x+i(A(y)-A(x))} h(y) dy,$$

where  $A(x) = \int_0^x a(t) dt$ . This is a version of  $\mathcal{H}_\Gamma$ ,  $\Gamma = \{x + iA(x); x \in \mathbb{R}\}$  by Calderón [C]. The following inequality [Mu3, p. 53] is established:  $\|C[a]\| \leq \text{Const}\{1 + \sqrt{\|a\|_{\text{BMO}}}\}$ . The proof in [Mu3] is not short, however, the method is founded on only one principle “the Calderón-Zygmund decomposition”. Using the separation theorem and the Calderón-Zygmund decomposition, we can deduce the required inequality. In order to see the exactness of the power  $3/2$ , it is sufficient to construct a compact set  $Q_n^*$  and a graph  $\Gamma_n^*$  containing  $Q_n^*$  so that  $\gamma(Q_n^*) \leq \text{Const}/\sqrt{n}$ ,  $|\text{pr } Q_n^*| = 1$  and  $|\Gamma_n^*| \leq \text{Const } n$ . Our example is related to David's example [D]: *For any  $M \geq 1$ , there exists a real-valued function  $a_M \in \text{BMO}$  such that  $\|a_M\|_{\text{BMO}} \leq M$  and  $\|C[a_M]\| \geq \text{Const}\sqrt{M}$ .* In the construction of  $(Q_n^*, \Gamma_n^*)$ , the following fact [Mu4] is important:  $\lim_{m \rightarrow \infty} \gamma(R(m)) < 1/4$ , where  $R(m) = \{x + i2^{-m}\text{sign}(\sin \pi 2^m x); 0 \leq x \leq 1\}^{\circ\ell}$ . (Note that  $\{R(m)\}_{m=1}^\infty$  converges to  $[0, 1]$  and  $\gamma([0, 1]) = 1/4$ .) Put

$$R(m_1, \dots, m_n) = \left\{ x + i \sum_{k=1}^n 2^{-m_1 - \dots - m_k} \text{sign}(\sin \pi 2^{m_1 + \dots + m_k} x); 0 \leq x \leq 1 \right\}^{\circ\ell}.$$

Choosing a sequence  $\{m_k^*\}_{k=1}^\infty$  of positive integers so that  $\{m_{k+1}^*/m_k^*\}_{k=1}^\infty$  is rapidly increasing, we put  $Q_n^* = R(m_1^*, \dots, m_n^*)$  ( $n \geq 1$ ). Then (3) yields that  $\gamma(Q_{n+1}^*) \leq \gamma(Q_n^*) - \delta_0 \gamma(Q_n^*)^3$  ( $n \geq 1$ ) with a small constant  $\delta_0$ . Thus  $\gamma(Q_n^*) \leq \text{Const}/\sqrt{n}$  ( $n \geq 1$ ). Connecting endpoints of  $Q_n^*$  by segments parallel to the  $y$ -axis, we obtain an arc  $\Gamma_n^*$  of length less than  $\text{Const. } n$ , which we can regard as a graph. Thus the power  $3/2$  is best possible. It is interesting to try to deduce the exactness of  $3/2$  by the dipole functions.

## 5. The arc-Length Variation

Let  $\mathcal{F}$  denote the totality of domains  $\Omega$  such that  $\Omega^c \in \mathcal{A}$ . For  $\Omega \in \mathcal{F}$ , let  $K(z, \bar{z}; \Omega)$  denote the reproducing kernel of  $H_0^2(\Omega) = \{f \in H^2(\Omega); f(\infty) = 0\}$  with respect to

$|dz|/(2\pi)$ , i.e.

$$f(\zeta) = \frac{1}{2\pi} \int_{\partial\Omega} \overline{K(z, \bar{\zeta}; \Omega)} f(z) |dz| \quad (f \in H_0^2(\Omega)), \quad K(\infty, \bar{\zeta}; \Omega) = 0.$$

This is called the Szegö kernel of  $H_0^2(\Omega)$  [B, Fa]. We have

$$\gamma(\Omega^c) = c(\infty; \Omega) = \frac{1}{(2\pi)^2} \int_{\partial\Omega} \int_{\partial\Omega} K(z, \bar{\zeta}; \Omega) dz d\bar{\zeta}.$$

Thus it is important to study the Szegö kernel of  $H_0^2(\Omega)$ . There are many articles about the variational approach to various Szegö kernels [GS, HS, S1, 2, Sm, SS]. We here show a variational formula for  $\gamma(E)$ ,  $E \in \mathcal{A}$  with respect to the arc-length. This is the variation of degenerate boundaries and related to Löwner's differential equation. There exists uniquely a pair  $(g(\cdot; \Omega), \phi(\cdot; \Omega))$  of functions in  $H^2(\Omega)$  such that  $g(\infty; \Omega) = 0$ ,  $\phi(\infty; \Omega) = 1$  and  $\frac{1}{i} \phi(z; \Omega) dz = \overline{g(z; \Omega)} |dz|$  a.e. on  $\partial\Omega$ . We have  $\psi(z; \Omega) = \phi(z; \Omega)^2$  and  $f(z; \Omega) = g(z; \Omega)/\phi(z; \Omega)$  [G]. For  $\zeta \in \Omega - \{\infty\}$ , there exists uniquely a pair  $(K(\cdot, \bar{\zeta}; \Omega), L(\cdot, \zeta; \Omega))$  of functions such that  $K(\cdot, \bar{\zeta}; \Omega) \in H_0^2(\Omega)$ ,  $(\cdot - \zeta)L(\cdot, \zeta; \Omega) \in H^2(\Omega)$ ,  $L(z, \zeta; \Omega) = \frac{1}{z - \zeta} + (\text{regular terms})$  near  $\zeta$  and  $\frac{1}{i} L(z, \zeta; \Omega) dz = \overline{K(z, \bar{\zeta}; \Omega)} |dz|$  a.e. on  $\partial\Omega$  [B]. The function  $K(z, \bar{\zeta}; \Omega)$  is none other than the Szegö kernel defined above. Here are two functions  $Dc$  and  $D^2c$  necessary for our variation. For three distinct numbers  $w, z, \zeta \in \Omega$ , we define

$$Dc(z, \zeta; \Omega) = |L(z, \zeta; \Omega)|^2 - |K(z, \bar{\zeta}; \Omega)|^2,$$

$$D^2c(w, z, \zeta; \Omega) = 2 \operatorname{Re} \{ DL(w, z, \zeta; \Omega) \overline{L(z, \zeta; \Omega)} - DK(w, z, \zeta; \Omega) \overline{K(z, \bar{\zeta}; \Omega)} \},$$

where

$$DK(w, z, \zeta; \Omega) = L(w, z; \Omega) \overline{L(w, \zeta; \Omega)} - \overline{K(w, \bar{z}; \Omega)} K(w, \bar{\zeta}; \Omega),$$

$$DL(w, z, \zeta; \Omega) = L(w, z; \Omega) \overline{K(w, \bar{\zeta}; \Omega)} - \overline{K(w, \bar{z}; \Omega)} L(w, \zeta; \Omega).$$

In the definition, we replace  $K(\cdot, \bar{\infty}; \Omega) = \overline{K(\infty, \cdot; \Omega)}$  by  $-g(\cdot; \Omega)$ , and replace  $L(\cdot, \infty; \Omega) = -L(\infty, \cdot; \Omega)$  by  $-\phi(\cdot; \Omega)$  if one of  $w, z, \zeta$  is  $\infty$ . Let  $\Gamma$  be a closed analytic arc such that  $\Gamma \subset \Omega^c$ ,  $\Omega - \Gamma \in \mathcal{F}$  and  $\Gamma \cap \Omega^c$  is at most a singleton. A continuous function  $w_t \in \Gamma$  on  $[0, |\Gamma|]$  is called the arc-length representation of  $\Gamma$  if  $w_0, w_{|\Gamma|}$  are endpoints of  $\Gamma$  and  $|\Gamma_t| = t$  ( $0 \leq t \leq |\Gamma|$ ), where  $\Gamma_t = \{w_s; 0 \leq s \leq t\}$ ; we define  $w_t$  so that  $\Gamma \cap \Omega^c = \{w_0\}$  if  $\Gamma \cap \Omega^c \neq \emptyset$ . We write  $\Omega_t = \Omega - \Gamma_t$  ( $0 \leq t \leq |\Gamma|$ ).

**Theorem 2** [Mu7].

(6) For any  $0 < t \leq |\Gamma|$ , the derivative  $\partial c(\infty; \Omega_t)/\partial t$ , the limit  $\lim_{u \downarrow t} Dc(w_u, \infty; \Omega_t)$  ( $= Dc(w_t, \infty; \Omega_t)$ , say) exist and  $\partial c(\infty; \Omega_t)/\partial t = Dc(w_t, \infty; \Omega_t)/4$ . The right-derivative  $\partial c(\infty; \Omega_0)/\partial t$  at  $t = 0$  exists and  $\partial c(\infty; \Omega_t)/\partial t$  is continuous on  $[0, |\Gamma|]$ .

(7) For any  $0 < t \leq |\Gamma|$  and  $z \in \Omega - (\Gamma \cup \{\infty\})$ , the derivative  $\partial Dc(z, \infty; \Omega_t)/\partial t$ , the limit  $\lim_{u \downarrow t} D^2c(w_u, z, \infty; \Omega_t)$  ( $= D^2c(w_t, z, \infty; \Omega_t)$ , say) exist and  $\partial Dc(z, \infty; \Omega_t)/\partial t =$

$D^2c(w_t, z, \infty; \Omega_t)/4$ . For any  $z \in \Omega - (\Gamma \cup \{\infty\})$ , the right-derivative  $\partial Dc(z, \infty; \Omega_t)/\partial t$  at  $t = 0$  exists and  $\partial Dc(z, \infty; \Omega_t)/\partial t$  is continuous on  $[0, |\Gamma|]$ .

In the case of  $t = 0$  and  $w_0 \in \Omega$ , the above formulae correspond to the variation by cutting a hole [HS, Sm, SS, p. 283]. Our method is based on the comparison (by the aid of conformal mappings) with the variation of segments. This theorem is applied as follows. Given  $E \in \mathcal{A}$ , we can write  $E = C_1 \cup \dots \cup C_n$  with mutually disjoint closed arcs  $\{C_j\}_{j=1}^n$ . Using the arc-length representations of  $C_j$  ( $j = 1, \dots, n$ ), we define a right-continuous arc-length representation  $W_t$  ( $0 \leq t \leq |E|$ ) of  $E$ . Then (6) shows that

$$(8) \quad \gamma(E) = c(\infty; E^c) = \frac{1}{4} \int_0^{|E|} Dc(W_t, \infty; E_t^c) dt,$$

where  $E_t = \{W_s; 0 \leq s \leq t\}$ . (The equality  $\gamma(F \cup \{w\}) = \gamma(F)$  is important in our argument.) Recall that  $Dc(W_t, \infty; E_t^c)$  is defined by the limit  $\lim_{u \downarrow t} Dc(W_u, \infty; E_t^c)$ . Given  $0 < t < |E|$  and  $z \in E_t^c - \{\infty\}$ , we take a right-continuous arc-length representation  $W_s^*$  ( $0 \leq s \leq t$ ) of  $E_t$ ;  $W_s^*$  may not be equal to  $W_s$ . Then (7) shows that

$$(9) \quad Dc(z, \infty; E_t^c) = Dc(z, \infty; F_s^c) + \frac{1}{4} \int_s^t D^2c(W_u^*, z, \infty; F_u^c) du,$$

where  $F_u = \{W_x^*; 0 \leq x \leq u\}$ . Using (8) and (9), we can study  $\gamma(E)$ . To investigate  $Dc$  and  $D^2c$ , we introduce a class  $\mathcal{G}$  of domains. Let  $\mathcal{G}$  denote the totality of domains  $\Omega$  with the following property;  $\Omega$  is expressed as  $\Omega = \Omega^* - E$  with  $E \in \mathcal{A}$  and a domain  $\Omega^* \supset E$  bounded by a finite number of Jordan curves  $\{C_j\}_{j=1}^n$  such that each  $C_j$  is a finite union of analytic arcs. The functions  $Dc$  and  $D^2c$  are defined for domains in  $\mathcal{G}$ . We see that  $Dc(z, \zeta; \Omega)|dz||d\zeta|$  and  $D^2c(w, z, \zeta; \Omega)|dw||dz||d\zeta|$  are conformally invariant. Thus we may discuss these differential forms in canonical domains. Now we show an application of Theorem 2. It is unknown whether  $\gamma(\cdot)$  is subadditive [Da, DØ, V2]. Saita [Su3] shows that  $\gamma(A \cup B) \leq \gamma(A) + \gamma(B)$  if  $A$  and  $B$  are disjoint continua. Equalities (8) and (9) yield that  $\gamma(\cdot)$  is subadditive if  $D^2c \leq 0$  for any domain.

**Theorem 3** [Mu7]. *The inequality  $D^2c \leq 0$  holds for simply and doubly connected domains.*

Applying this theorem to simply-connected domains, we see that  $\text{Cap}(A \cup B) \leq \text{Cap}(A) + \text{Cap}(B)$  if  $A$  and  $B$  are two continua with an intersection, where  $\text{Cap}(\cdot)$  is logarithmic capacity [Z, p. 134]. Note that  $\text{Cap}(\cdot)$  is not subadditive. The case of doubly-connected domains shows Saita's subadditivity and yields that  $\gamma(A \cup B) \leq \gamma(A) + \gamma(B)$  if  $A$  is a union of two continua and  $B$  is a continuum intersecting with  $A$ . If  $\Omega$  is simply-connected, then we may assume that  $\Omega$  is the open unit disk  $\mathbb{D}$  and  $w = 0$ . We have

$$\begin{aligned} D^2c(0, z, \zeta; \mathbb{D}) &= - \frac{(1 + |z\zeta|)^4 |1 + z\bar{\zeta}|^2}{(1 - |z|^2)(1 - |\zeta|^2)|z\zeta|^2|z - \zeta|^2} \\ &\times \left\{ \left( \frac{|z| + |\zeta|}{1 + |z\zeta|} \right)^2 - \left| \frac{z + \zeta}{1 + z\bar{\zeta}} \right|^2 \right\} \left\{ \left( \frac{|z| + |\zeta|}{1 + |z\zeta|} \right)^2 - \left| \frac{z - \zeta}{1 - z\bar{\zeta}} \right|^2 \right\} \quad (\leq 0). \end{aligned}$$

The hyperbolic distance in  $\mathbb{ID}$  is defined by  $d(z, \zeta) = \operatorname{arctanh}(|z - \zeta|/|1 - z\bar{\zeta}|)$  ( $z, \zeta \in \mathbb{ID}$ ). We have  $d(z, 0) + d(0, \zeta) = \operatorname{arctanh}\{(|z| + |\zeta|)/(1 + |z\zeta|)\}$ . Thus the inequality  $D^2c \leq 0$  is related to the triangle inequality with respect to  $d(\cdot, \cdot)$ . If  $\Omega$  is doubly-connected, then we may assume that  $\Omega$  is a ring  $R_\varrho = \{\varrho < |z| < 1\}$ . We have

$$\begin{aligned} D^2c(w, z, \zeta; R_\varrho) &= \frac{2K^3}{\pi^3 |wz\zeta|} \\ &\times \operatorname{Im} \left[ \left( \frac{dn(\xi - u)}{sn(\xi - u)} \frac{dn(\bar{\xi} - v)}{sn(\bar{\xi} - v)} - \frac{dn(\bar{\xi} - u)}{sn(\bar{\xi} - u)} \frac{dn(\xi - v)}{sn(\xi - v)} \right) \frac{dn(\bar{u} - \bar{v})}{sn(\bar{u} - \bar{v})} \right. \\ &\quad \left. - \left( \frac{dn(\xi - u)}{sn(\xi - u)} \frac{dn(\bar{\xi} - \bar{v})}{sn(\bar{\xi} - \bar{v})} - \frac{dn(\bar{\xi} - u)}{sn(\bar{\xi} - u)} \frac{dn(\xi - \bar{v})}{sn(\xi - \bar{v})} \right) \frac{dn(\bar{u} - v)}{sn(\bar{u} - v)} \right], \\ \xi &= \frac{K}{i\pi} \log w, \quad u = \frac{K}{i\pi} \log z, \quad v = \frac{K}{i\pi} \log \zeta, \quad K = K(k), \end{aligned}$$

where the modulus  $k$  is defined by  $\log \varrho = -\pi K(\sqrt{1 - k^2})/K(k)$ . The following expression is also applicable to compute  $D^2c$ : Let  $\Omega^* = \{\bigcup_{k=1}^n [a_k, b_k]\}^c$  ( $a_1 < b_1 < \dots < a_n < b_n$ ). Then

$$\begin{aligned} (10) \quad D^2c(w, z, \zeta; \Omega^*) &= \frac{1}{4|M(w)M(z)M(\zeta)|} \\ &\times \operatorname{Re} \left[ \left( \frac{M(w) + M(z)}{w - z} \frac{\overline{M(w)} - M(\zeta)}{\bar{w} - \zeta} \right. \right. \\ &\quad \left. \left. - \frac{\overline{M(w)} - M(z)}{\bar{w} - z} \frac{M(w) + M(\zeta)}{w - \zeta} \right) \frac{\overline{M(z)} + \overline{M(\zeta)}}{\bar{z} - \bar{\zeta}} \right. \\ &\quad \left. - \left( \frac{M(w) + M(z)}{w - z} \frac{\overline{M(w)} + \overline{M(\zeta)}}{\bar{w} - \bar{\zeta}} \right. \right. \\ &\quad \left. \left. - \frac{\overline{M(w)} - M(z)}{\bar{w} - z} \frac{M(w) - \overline{M(\zeta)}}{w - \bar{\zeta}} \right) \frac{\overline{M(z)} - M(\zeta)}{\bar{z} - \zeta} \right], \\ M(\zeta) &= \prod_{k=1}^n \sqrt{(b_k - \zeta)/(a_k - \zeta)} \quad (\zeta = w, z, \zeta). \end{aligned}$$

The ring  $R_\varrho$  is conformally mapped to a domain of this type with  $n = 2$ . To see  $D^2c(w, z, \zeta; \Omega^*) \leq 0$  in the case of  $n = 2$ , we may assume that  $w = x \in \mathbb{R}$ ,  $z/i > 0$  and  $\zeta/i < 0$ . Then (10) shows that

$$D^2c(x, z, \zeta; \Omega^*) = \frac{Ax^2 + Bx + C}{|M(z)M(\zeta)||(x - z)(x - \zeta)|^2} \quad (A, B, C \in \mathbb{R}).$$

By an inequality of Möbius type in elementary geometry, we obtain  $A < 0$  and  $B^2 - 4AC < 0$ .

## References

- [A] L. Ahlfors: Bounded analytic functions. Duke. Math. J. **14** (1947) 1–11
- [AB] L. Ahlfors, A. Beurling: Conformal invariants and function-theoretic null-sets. Acta Math. **83** (1950) 101–129
- [B] S. Bergman: The kernel function and conformal mapping. Math. Surv. V. Amer. Math. Soc., New York 1950
- [C] A.P. Calderón: Commutators, singular integrals on Lipschitz curves and applications. ICM Helsinki 1978, pp. 85–96
- [Ch] M. Christ: Lectures on singular integral operators. CBMS 77. Amer. Math. Soc., Providence 1990
- [CJS] R.R. Coifman, P.W. Jones, S. Semmes: Two elementary proofs of the  $L^2$  boundedness of Cauchy integrals on Lipschitz curves. J. Amer. Math. Soc. **2** (1989) 553–564
- [CMM] R.R. Coifman, A. McIntosh, Y. Meyer: L'intégrale de Cauchy définit un opérateur borné sur  $L^2$  pour les courbes lipschitziennes. Ann. Math. **116** (1982) 361–388
- [D] G. David: Opérateurs de Calderón-Zygmund. ICM Berkeley 1986, pp. 890–899
- [Da] A.M. Davie: Analytic capacity and approximation problems. Trans. Amer. Math. Soc. **171** (1972) 409–444
- [DØ] A.M. Davie, B. Øksendal: Analytic capacity and differentiability properties of finely harmonic functions. Acta Math. **149** (1982) 127–152
- [F] K.J. Falconer: The geometry of fractal sets. Cambridge Univ. Press, Cambridge 1985
- [Fa] J.D. Fay: Theta functions on Riemann surfaces. (Lecture Notes in Mathematics, vol. 352.) Springer, Berlin Heidelberg New York 1973
- [Fe] C. Ferrari: Sulla trasformazione conforme di due cerchi in due profili alari. Memorie della R. Accad. delle Scienze di Torino Serie II **67** (1930) 1–15
- [G] P. Garabedian: Schwarz's lemma and the Szegő kernel function. Trans. Amer. Math. Soc. **67** (1949) 1–35
- [Ga] J. Garnett: Analytic capacity and measure. (Lecture Notes in Mathematics, vol. 297.) Springer, Berlin Heidelberg New York 1972
- [Gar] I.E. Garrick: Potential flow about arbitrary biplane wing sections. Technical Report No. 542, NACA 1936, pp. 47–75
- [GS] P. Garabedian, M. Schiffer: Identities in the theory of conformal mapping. Trans. Amer. Math. Soc. **65** (1949) 187–238
- [H] S.Ya. Havinson: Analytic capacity of sets, joint nontriviality of various classes of analytic functions and the Schwarz lemma in arbitrary domains (Russian) Mat. Sb. **54** (1961) 3–50
- [HS] N.S. Hawley, M. Schiffer: Half-order differentials on Riemann surfaces. Acta Math. **115** (1966) 199–236
- [HHN] V.P. Havin, S.V. Hruščëv, N.K. Nikol'skii: Linear and complex analysis problem book. (Lecture Notes in Mathematics, vol. 1043.) Springer, Berlin Heidelberg New York 1984
- [J] P.W. Jones: Square functions, Cauchy integrals, analytic capacity and harmonic measure. In: Harmonic analysis and partial differential equations, pp. 24–68 (Lecture Notes in Mathematics, vol. 1384.) Springer, Berlin Heidelberg New York 1989
- [JM] P.W. Jones, T. Murai: Positive analytic capacity but zero Buffon needle probability. Pacific J. Math. **133** (1988) 99–114
- [M1] P. Mattila: A class of sets with positive length and zero analytic capacity. Ann. Acad. Sci. Fenn. Ser. AI **10** (1985) 387–395
- [M2] P. Mattila: Smooth maps, null-sets for integralgeometric measure and analytic capacity. Ann. Math. **123** (1986) 303–309

- [Ma] D.E. Marshall: Removable sets for bounded analytic functions. In: Linear and complex analysis problem book, pp. 485–490 (Lecture Notes in Mathematics, vol. 1043.) Springer, Berlin Heidelberg New York 1984
- [Me] Y. Meyer: Wavelets and operators. In: Analysis at Urbana 1 pp. 256–365 (Lecture Note Ser. 137.) London Math. Soc., Cambridge 1989
- [Mi] L.M. Milne-Thomson: Theoretical hydrodynamics, Second edition. Macmillan, London 1949
- [Mu1] T. Murai: Construction of  $H^1$  functions concerning the estimate of analytic capacity. Bull. London Math. Soc. **19** (1987) 154–160
- [Mu2] T. Murai: Comparison between analytic capacity and the Buffon needle probability. Trans. Amer. Math. Soc. **304** (1987) 501–514
- [Mu3] T. Murai: A real variable method for the Cauchy transform, and analytic capacity. (Lecture Notes in Mathematics, vol. 1307.) Springer, Berlin Heidelberg New York 1988
- [Mu4] T. Murai: The power 3/2 appearing in the estimate of analytic capacity. Pacific J. Math. **143** (1990) 313–340
- [Mu5] T. Murai: Analytic capacity for two segments. Nagoya Math. J. **122** (1991), to appear
- [Mu6] T. Murai: A formula for analytic separation capacity. Kōdai Math. J. **13** (1990) 265–288
- [Mu7] T. Murai: The arc-length variation of analytic capacity and a conformal geometry. Submitted
- [P] Ch. Pommerenke: Über die analytische Kapazität. Arch. Math. **11** (1960) 270–277
- [S1] M. Schiffer: Variational methods in the theory of Riemann surfaces. In: Contributions to the theory of Riemann surfaces, pp. 15–30. (Ann. Math. Stud. 30.) Princeton Univ. Press, Princeton 1953
- [S2] M. Schiffer: Some recent developments in the theory of conformal mapping. In: R. Courant, Dirichlet's principle, pp. 249–318. (Pure Appl. Math. III) Interscience, New York 1967
- [Sm] E.P. Smith: The Garabedian function of an arbitrary compact set. Pacific J. Math. **51** (1974) 289–300
- [Su1] N. Suita: On a metric induced by analytic capacity. Kōdai Math. Sem. Rep. **25** (1973) 215–218
- [Su2] N. Suita: On a metric induced by analytic capacity II. Kōdai Math. Sem. Rep. **27** (1976) 159–162
- [Su3] N. Suita: On subadditivity of analytic capacity for two continua. Kōdai Math. J. **7** (1984) 73–75
- [SO] L. Sario, K. Oikawa: Capacity functions. Springer, Berlin Heidelberg New York 1969
- [SS] M. Schiffer, D.C. Spencer, Functionals of finite Riemann surfaces. Princeton Univ. Press, Princeton 1954
- [V1] A.G. Vitushkin: Example of a set of positive length but of zero analytic capacity (Russian). Dokl. Akad. Nauk SSSR **127** (1959) 246–249
- [V2] A.G. Vitushkin: Analytic capacity of sets in problems of approximation theory (Russian). Uspehi Mat. Nauk **22** (1967) 141–199
- [Z] L. Zalcman: Analytic capacity and rational approximation. (Lecture Notes in Mathematics, vol. 50.) Springer, Berlin Heidelberg New York 1968



# Recent Applications of $L^2$ Estimates for the Operator $\bar{\partial}$

Takeo Ohsawa

Department of Mathematics, Nagoya University, Nagoya 464-01, Japan

1. In the theory of holomorphic functions of several variables, boundaries of complex manifolds arose as the singularities of analytic objects. Geometric structure of manifolds with boundaries of this sort is therefore of interest in function theory. As for the singularity of holomorphic functions, a general picture was given by the solution of the Levi problem given by Oka [39] over  $C^n$  and by Grauert [16] on complex manifolds, which characterized Stein manifolds by the existence of strictly plurisubharmonic exhaustion functions. The latter's work was based on the cohomology finiteness theorem and was generalized by Andreotti-Grauert [2] to noncompact complex spaces that admit certain exhaustion functions. As is well known, a method of partial differential equation is available to study the cohomology groups of Riemannian manifolds (cf. [7]). Importance of this method in function theory became apparent by the works of Andreotti-Vesentini [4] and Hörmander [19]. The latter is already a penetrating work that recovers main results of [39, 16] and [2] in a completely different way. The method consists in establishing an à priori  $L^2$  estimate for a given  $\bar{\partial}$ -closed form, which is usually a direct consequence of commutator relations in a graded operator algebra generated by several covariant exterior differential operators like  $\bar{\partial}$ ,  $\partial$  and their adjoints. By and by it has turned out that this approach, namely the  $L^2$  theory, has an advantage in obtaining more detailed information about the analytic cohomology groups and functions on noncompact complex manifolds (cf. [21, 25, 17, and 43]). The basic problems in this context are therefore to clarify the specific properties of the  $L^2$  objects, like  $L^2$  cohomology groups and harmonic forms, and relationship between the  $L^2$  and the ordinary cohomology groups. We shall report below on recent results about the  $L^2$  cohomology groups of noncompact manifolds with emphasis on the extension of the classical Hodge theory which turned out to have an application to intersection cohomology theory.
2. Basic technical devices are summarized here. Let  $(X, ds^2)$  be a connected Hermitian manifold of dimension  $n$ , and let  $(L, h)$  be a Hermitian line bundle over  $X$ . For any square integrable  $L$ -valued  $(p, q)$ -form  $u$ ,  $\|u\|_h$  will denote the  $L^2$  norm of  $u$  with respect to  $ds^2$  and  $h$ . The fundamental form of  $ds^2$  will be denoted by  $\omega$  and the curvature form of  $h$  by  $\Theta_h$ .

**Theorem 1.** If  $ds^2$  is a complete Kähler metric on  $X$  and  $i\Theta_h = \omega$ , then for any  $L$ -valued (resp.  $L^{-1}$ -valued)  $(p, q)$ -form  $v$  with  $\bar{\partial}v = 0$   $\deg v$  ( $:= p + q$ )  $> n$  (resp.  $\deg v < n$ ) and  $\|v\|_h < \infty$ , there exists an  $L$ -valued (resp.  $L^{-1}$ -valued)  $(p, q - 1)$ -form  $u$  satisfying  $\bar{\partial}u = v$  and  $\|u\|_h \leq \|v\|_h$ . Moreover the range of the maximal closed extension of  $\bar{\partial}|C_0(X, L^{-1})$  is closed in the degrees  $\leq n$ .

The proof of Theorem 1 is based on the complex Weitzenböck's formula which we recall briefly. Let  $C_0(X, L)$  denote the set of  $L$ -valued  $C^\infty$  differential forms with compact support on  $X$ , let  $\bar{\partial}_h^*$  denote the Hilbert space adjoint of  $\bar{\partial}|C_0(X, L)$  with respect to  $ds^2$  and  $h$ , and let  $\Lambda$  be the (pointwise) adjoint of exterior multiplication by the fundamental form of  $ds^2$ . Then we have

$$[\bar{\partial}, \bar{\partial}_h^*]_{\text{gr}} - *^{-1}[\bar{\partial}, \bar{\partial}_h^*]_{\text{gr}}* = [i\Theta_h, \Lambda]_{\text{gr}} \quad (1)$$

on  $C_0(X, L)$ , where  $[a, b] := ab - (-1)^{\deg a \deg b}ba$ ,  $*$  denotes the Hodge's star operator and  $i\Theta_h$  is identified with the corresponding exterior multiplication. The à priori estimates needed for the above existence theorems follow directly from (1). Originally, Theorem 1 was stated only for compact manifolds for it was thought of as an analytic counterpart of Lefschetz's hyperplane section theorem on nonsingular projective varieties (cf. Akizuki-Nakano [1]). The above noncompact version is due to Andreotti-Vesentini [3] and recently it was used in an essential way in the proof of Cheeger-Goresky-MacPherson's conjecture for varieties with isolated singularities, as we shall see later. The following looks very likely to be a corollary of Theorem 1, but it was discovered much later by Donnelly-Fefferman [14] in the study of Schwartz kernels on strongly pseudoconvex domains, and the proof is actually independent of Theorem 1.

**Theorem 2.** Under the situation of Theorem 1, assume particularly that  $L$  is the trivial bundle so that the connection form  $\partial \log h$  is identified with a  $(1, 0)$ -form on  $X$ , and that  $|\partial \log h|$  is bounded. Then for any  $\bar{\partial}$ -closed  $(p, q)$ -form  $v$  with  $p + q \neq n$  which is square integrable with respect to the trivial fiber metric  $h_0$ , there exists a  $(p, q - 1)$ -form  $u$  satisfying  $\bar{\partial}u = v$  and  $\|u\|_{h_0} \leq \|v\|_{h_0}$ . Moreover the range of the maximal closed extension of  $\bar{\partial}$  with respect to the metric  $h_0$  is closed.

The proof of Theorem 2 is immediate from the equality

$$[\bar{\partial}, (\bar{\partial} \log h)^*]_{\text{gr}} + *^{-1}[\bar{\partial}, (\bar{\partial} \log h)^*]_{\text{gr}}* = [i\Theta_h, \Lambda]_{\text{gr}}. \quad (2)$$

Another important machinery in the  $L^2$  theory is the following.

**Theorem 3.** Assume that  $i\Theta_h = \omega$  for given Hermitian metrics  $ds^2$  and  $h$ . If  $X$  admits a complete Kähler metric, then for any  $L$ -valued  $(n, q)$ -form  $v$  with  $\bar{\partial}v = 0$ ,  $q \geq 1$  and  $\|v\|_h < \infty$  (with respect to  $ds^2$  and  $h$ ) there exists an  $L$ -valued  $(n, q - 1)$ -form  $u$  satisfying  $\bar{\partial}u = v$  and  $\|u\|_h \leq \|v\|_h$ .

Theorem 3 contains Hörmander's theorem (cf. [19]) whose applications are already widespread. The above formulation is given by Demainay [11] and Ohsawa [27, 31] independently, aiming at applying Hörmander's method to more geometric

questions related to holomorphic  $n$ -forms and Bergman kernels (see also [30, 13, 15]). It is also regarded as a noncompact version of Kodaira vanishing theorem.

**3.** We are going to sketch how to apply the above mentioned tools to relate the ordinary and  $L^2$  cohomology groups of Hermitian manifolds  $(X, ds^2)$ . Let  $C_0(X)$  be the set of  $C^\infty$  forms on  $X$  with compact support and let  $\bar{\partial}_{\max}$  be the maximal closed extension of  $\bar{\partial}|_{C_0(X)}$  to the space of square integrable forms  $L_{(2)}(X)$ . Namely the domain of  $\bar{\partial}_{\max}$ , denoted by  $\text{Dom } \bar{\partial}_{\max}$ , consists of square integrable forms  $u$  for which  $\bar{\partial}u$  is also square integrable. The  $(p, q)$ -component of  $L_{(2)}(X)$  will be denoted by  $L^{p,q}(X)$ , and we call the space  $\text{Ker } \bar{\partial}_{\max} \cap L^{p,q}(X)/\text{Im } \bar{\partial}_{\max} \cap L^{p,q}(X)$  the  $L^2$   $\bar{\partial}$ -cohomology group of type  $(p, q)$ , denoted by  $H_{(2)}^{p,q}(X)$  or  $H_{(2)}^{p,q}(X)_{ds^2}$ . For any compact subset  $K \subset X$  we set

$$L^{p,q}(X/K) := \{u \in L^{p,q}(X); \text{supp } u \subset X \setminus K\}$$

and  $H_{(2)}^{p,q}(X/K) := \text{Ker } \bar{\partial}_{\max} \cap L^{p,q}(X/K)/\bar{\partial}(\text{Dom } \bar{\partial}_{\max} \cap L^{p,q-1}(X))$ . Then we have a long exact sequence

$$\cdots \rightarrow \varprojlim_{K \Subset X} H_{(2)}^{p,q}(X/K) \rightarrow H_{(2)}^{p,q}(X) \rightarrow H^{p,q}(X) \rightarrow \varinjlim_{K \Subset X} H_{(2)}^{p,q+1}(X/K) \rightarrow \cdots$$

where  $H^{p,q}(X)$  denotes the ordinary  $\bar{\partial}$ -cohomology group of type  $(p, q)$ . As for the  $\bar{\partial}$ -cohomology groups with compact support  $H_0^{p,q}(X)$ , we have

$$\cdots \rightarrow \varinjlim_{K \Subset X} H_{(2)}^{p,q-1}(X \setminus K) \rightarrow H_0^{p,q}(X) \rightarrow H_{(2)}^{p,q}(X) \rightarrow \varprojlim_{K \Subset X} H_{(2)}^{p,q}(X \setminus K) \rightarrow \cdots$$

Hence we have the following criterion.

**Proposition 4.** *The canonical homomorphism  $\alpha^{p,q}$  (resp.  $\beta^{p,q}$ ) is surjective if  $\lim_{\leftarrow} H_{(2)}^{p,q+1}(X/K) = 0$  (resp. if  $\lim_{\rightarrow} H_{(2)}^{p,q}(X \setminus K) = 0$ ) and injective if  $\lim_{\leftarrow} H_{(2)}^{p,q}(X/K) = 0$  (resp. if  $\lim_{\rightarrow} H_{(2)}^{p,q-1}(X \setminus K) = 0$ ).*

In case the boundary of  $X$  is “small”, conditions of Proposition 4 are actually satisfied.

**Proposition 5.** *Suppose that  $X$  is the regular part of a projective variety  $Z$  such that  $\dim(Z \setminus X) = 0$ . Then there exists a complete Kähler metric on  $X$  for which  $\lim_{\leftarrow} H_{(2)}^{p,q}(X/K) = 0$  if  $p + q < n$  and  $\lim_{\rightarrow} H_{(2)}^{p,q}(X \setminus K) = 0$  if  $p + q > n$ .*

We note that the à priori estimates for the  $(p, q)$ -forms on  $X \setminus K$  imply the finite dimensionality of  $H_{(2)}^{p,q}(X)$  and the separatedness of  $H_{(2)}^{p,q+1}(X)$ , so that we have the following straightforward consequence of Proposition 5.

**Theorem 6** (cf [32]). *Let  $X$  be as in Proposition 5. Then there exists a complete Kähler metric on  $X$  such that*

- 1) *The canonical homomorphisms  $\alpha^{p,q}$  (resp.  $\beta^{p,q}$ ) are bijective if  $p + q < n - 1$  (resp. if  $p + q > n + 1$ ) and injective (resp. surjective) if  $p + q = n - 1$  (resp.  $p + q = n + 1$ ).*

2)  $H_{(2)}^{p,q}(X)$  are finite dimensional if  $p + q \neq n$  and separated if  $p + q = n$ .

The above metric is constructed as the sum of the Fubini-Study metric (restricted to  $X$ ) and the complex Hessian of a  $C^\infty$  exhaustion function with values in  $(-\infty, 0]$ , say  $\phi$  which behaves like  $-\log(-\log \delta)$  near the points of  $Z \setminus X$ , where  $\delta$  denotes the distance to  $Z \setminus X$ . Theorem 2 is then applied for a sufficiently small sublevel set of  $\phi$  equipped with a complete Kähler metric of the form  $\partial\bar{\partial}\lambda(\phi)$  for a suitable convex increasing function  $\lambda$ .

By weakening the assumptions of Theorems 1 and 2 one has a more general  $L^2$  vanishing theorem that yields the following in a similar manner.

**Theorem 7** (cf. [33]). *If  $X$  is the regular part of a projective variety  $Z$  such that  $\dim(Z \setminus X) \leq k$ , there exists a complete Kähler metric on  $X$  such that*

1)  $\alpha^{p,q}$  (resp.  $\beta^{p,q}$ ) are bijective if  $p + q < n - k - 1$  (resp.  $p + q > n + k + 1$ ) and injective (resp. surjective) if  $p + q = n - k - 1$  (resp.  $p + q = n + k + 1$ )

2)  $H_{(2)}^{p,q}(X)$  are finite dimensional if  $|p + q - n| > k + 1$  and separated if  $p + q = n - k - 1$ .

Let  $H^r(X)$  and  $H_0^r(X)$  denote respectively the  $r$ -th cohomology group of  $X$  and that with compact support. Then Theorem 7 implies an extension of Hodge theory to quasiprojective varieties. Namely, as a corollary of Theorem 7 we obtain

**Theorem 8.** *Under the situation of Theorem 7,*

1)  $H^{p,q}(X) \cong H^{q,p}(X)$  for  $p + q < n - k - 1$  and  $H_0^{p,q}(X) \cong \overline{H_0^{q,p}(X)}$  for  $p + q > n + k + 1$ .

2)  $H^r(X)$  (resp.  $H_0^r(X)$ ) is canonically isomorphic to  $\bigoplus_{p+q=r} H^{p,q}(X)$  for  $r < n - k - 1$  (resp. isomorphic to  $\bigoplus_{p+q=r} H_0^{p,q}(X)$  for  $r > n + k + 1$ ).

**Remark.** There are different proofs of Theorem 8. Namely Bauer and Kosarew [6] use characteristic  $p$  method and Arapura [5] the technique of logarithmic differential forms.

4. The preceding discussion says nothing about the properties of  $H_{(2)}^{p,q}(X)$  for  $|p + q - n| \leq k$ . Since there is little hope to get any simple relation between the  $L^2$  and the ordinary cohomology in this range, we must compromise at the moment to study the  $L^2$  cohomology groups of  $(X, ds^2)$  for the exterior differential  $d$  instead of  $\bar{\partial}$ , which shall be denoted by  $H_{(2)}^r(X)$  or  $H_{(2)}^r(X)_{ds^2}$ . Cheeger-Goreski-MacPherson [10] posed the following fundamental question for the  $L^2$  cohomology with respect to the Fubini-Study metric. To fix our notation we set  $H_{(2)}^r(U) := H_{(2)}^r(X \cap U)_{ds^2}$  for open subsets  $U$  of  $Z$  if  $ds^2$  is the Fubini-Study metric.

**Cheeger-Goresky-MacPherson's Conjecture.** *Let  $Z \subset P^N$  be any irreducible projective variety and let  $X \subset Z$  be the set of regular points. Then  $H_{(2)}^r(Z)$  are canonically isomorphic to the intersection cohomology groups  $IH^r(Z)$  of the middle perversity (for the definition of  $IH^r(Z)$ , see [10]).*

In case  $\dim(Z \setminus X) = 0$ , the conjecture is equivalent to saying that

$$H_{(2)}^r(Z) \cong \begin{cases} H^r(X) & \text{if } r < n \\ \text{Im}(H_0^n(X) \rightarrow H^n(X)) & \text{if } r = n \\ H_0^r(X) & \text{if } r > n, \end{cases} \quad (3)$$

which has been verified in this case by Cheeger [9] Hsiang-Pati [18] and Nagase [24] for  $\dim Z \leq 2$  and by the author in April of 1990, and we are going to sketch the idea of the proof. As we have seen in the case of the  $L^2$   $\bar{\partial}$ -cohomology, (3) is equivalent to saying that  $\lim_{\leftarrow} H_{(2)}^r(Z/K) = 0$  for  $r < n$  and  $\lim_{\rightarrow} H_{(2)}^r(Z \setminus K) = 0$  for  $r \geq n$ . It is more or less routine that these vanishing follow from Theorem 1 except for the case  $r = n$  (cf. [32, 35]). Direct approaches to prove  $\lim H_{(2)}^r(Z \setminus K) = 0$  require subtle analysis as in [18] and [24]. Instead we take an indirect way. Namely we aim at proving the following.

**Proposition 9.** *There exists a Hermitian metric  $d\sigma^2$  on  $X$  such that  $\dim H_{(2)}^r(X)_{d\sigma^2} = \dim H_{(2)}^r(Z)$  and  $\dim H_{(2)}^r(X)_{d\sigma^2} = \dim IH^r(Z)$  for  $r = n \pm 1, n$ .*

A candidate of  $d\sigma^2$  is given by a result of Saper [41] (see also [23]):

**Theorem 10.** *If  $\dim(Z \setminus X) = 0$ , there exists a  $C^\infty$  exhaustion function  $\phi : X \rightarrow (-\infty, 0]$  such that  $d\sigma^2 := ds^2 + \bar{\partial}\partial\phi$  is a complete Kähler metric on  $X$  satisfying*

- 1) *The length of  $\partial\phi$  is bounded*
- 2)  *$H_{(2)}^r(X)_{d\sigma^2} \cong IH^r(Z)$ .*

In order to establish the required equalities for the  $L^2$  cohomology groups, we compare the spaces of harmonic forms by applying the approximation method of Runge-Hörmander. Namely let  $d\sigma_\varepsilon^2 = ds^2 + \varepsilon\partial\bar{\partial}\phi$  for  $\varepsilon \geq 0$  and let  $\mathcal{H}_\varepsilon^r$  be the set of harmonic  $r$ -forms with respect to  $d\sigma_\varepsilon^2$ . Then  $\mathcal{H}_\varepsilon^r \cong \mathcal{H}_1^r$  for  $\varepsilon > 0$  so that it suffices to show that  $\dim \mathcal{H}_0^r = \dim \mathcal{H}_1^r$ , for the range of  $d_{\max}$  (w.r.t.  $ds^2$ ) is closed (cf. [35]). Let  $\| \cdot \|_\varepsilon$  and  $d_\varepsilon^*$  denote respectively the  $L^2$  norm and the adjoint of  $d$  with respect to  $d\sigma_\varepsilon^2$ . Then the approximation argument proceeds as follows.

**Lemma.** *There exists a family of positive functions  $\Phi_\varepsilon$  on  $X$  for  $\varepsilon \searrow 0$  such that*

- 1)  $\|\Phi_\varepsilon u\|_\varepsilon \leq \text{const.}(\|u\|_\varepsilon + \|du\|_\varepsilon + \|d_\varepsilon^* u\|_\varepsilon)$  uniformly in  $\varepsilon$  for all  $u \in C_0(X)$  with  $\deg u \neq n$ .
- 2) *From any sequence of locally square integrable forms  $f_\varepsilon$  with  $\deg f_\varepsilon \neq n$  satisfying  $\|\Phi_\varepsilon f_\varepsilon\|_\varepsilon = 1$  and  $\|(d + d_\varepsilon^*)f_\varepsilon\|_\varepsilon < \infty$  one can choose a subsequence  $f_{\varepsilon_i}$  such that  $\{f_{\varepsilon_i}\}$  converges to a nonzero element of  $\text{Dom}(d_{\max} + d_{\max}^*)$  strongly on compact subsets of  $X$ .*

As  $\Phi_\varepsilon$ , we may take the length of  $\varepsilon\partial\phi - \partial \log(\log(\delta^{-1} + 1))$  with respect to  $d\sigma_\varepsilon^2$ . Here  $\delta$  is as before. The required property of  $\Phi_\varepsilon$  follows from the non-integrability of  $\delta^{-1} \log \delta$  on  $(0, 1/2)$  and an  $L^2$  estimate

$$\|\Phi_\varepsilon u\|_\varepsilon \leq \text{const.}(\|\chi_K u\|_\varepsilon + \|du\|_\varepsilon + \|d_\varepsilon^* u\|_\varepsilon),$$

which is valid uniformly in  $\varepsilon \searrow 0$  for some compact set  $K \subset X$  and all  $u \in C_0(X)$  with  $\deg u \neq n$ , where  $\chi_K$  stands for the characteristic function of  $K$ .

Suppose that  $\dim \mathcal{H}_1^{n+1} > \dim \mathcal{H}_0^{n+1}$ . Let  $U \subset X$  be any relatively compact open subset which contains  $K$ . Then, by hypothesis there must exist  $f_\varepsilon \in \mathcal{H}_\varepsilon^{n+1}$  for  $\varepsilon > 0$  such that

- 1)  $(\chi_U f_\varepsilon, g)_\varepsilon = 0$  for all  $g \in \mathcal{H}_0^{n+1}$
- 2)  $\|\Phi_\varepsilon f_\varepsilon\|_\varepsilon = 1$ .

Here  $(\ , )_\varepsilon$  denotes the inner product associated with  $\| \ \|_\varepsilon$ . But Lemma says then there exist  $f \in \mathcal{H}_0^{n+1}$  with  $f \neq 0$  satisfying  $(\chi_U f, f)_0 = 0$ , which contradicts the unique continuation theorem since  $f$  is harmonic. Since  $\dim H_0^{n+1}(X) \geq \dim \mathcal{H}_0^{n+1}$ , we may conclude that  $\dim \mathcal{H}_1^{n+1} = \dim \mathcal{H}_0^{n+1}$  in virtue of Saper's theorem. The proof for the case  $r = n - 1$  is similar. If  $r = n$ , suppose that  $\dim \mathcal{H}_0^n > \dim \mathcal{H}_1^n$ . Then there exists a finite dimensional subspace  $\mathcal{H}' \subset \mathcal{H}_0^n$  from which we can choose  $g_\varepsilon$  ( $\varepsilon \searrow 0$ ) satisfying  $\|g_\varepsilon\|_\varepsilon = 1$  and  $(\chi_U g_\varepsilon, g')_\varepsilon = 0$  for any  $g' \in \mathcal{H}_\varepsilon^n$ . Since  $\mathcal{H}_\varepsilon^{n+1}$  approximate  $\mathcal{H}_0^{n+1}$  in the above sense, it follows from (4) that there exist  $(u_\varepsilon, v_\varepsilon)$  satisfying  $du_\varepsilon + d_\varepsilon^* v_\varepsilon = \chi_U g_\varepsilon$  such that  $\|\Phi_\varepsilon u_\varepsilon\|_\varepsilon$  and  $\|\Phi_\varepsilon v_\varepsilon\|_\varepsilon$  are uniformly bounded in  $\varepsilon$ . Taking Lemma into account, this implies the existence of  $g \neq 0$  in  $\mathcal{H}'$  such that  $\chi_U g$  belongs to the range of  $d_{\max} + d_{\max}^*$ , which is an absurdity. Thus we get  $\dim \mathcal{H}_0^n \leq \dim \mathcal{H}_1^n$ . Clearly  $\dim \mathcal{H}_0^n \geq \dim \text{Im}(H_0^n(X) \rightarrow H^n(X))$ . Therefore one must have  $\dim \mathcal{H}_0^n = \dim \mathcal{H}_1^n$  by Saper's theorem.

Concerning the  $L^2$  cohomology with respect to the Fubini-Study metric, a more advanced question is whether there also exist geometric interpretation of the  $L^2$   $\bar{\partial}$ -cohomology groups  $H_{(2)}^{p,q}(Z)$  ( $:= H_{(2)}^{p,q}(X, ds^2)$ ). Recently Pardon-Stern [40] obtained the following.

**Theorem 11.** *Let  $Z$  be an irreducible projective variety of dimension  $n$  with isolated singularities. Then*

$$\chi''_{(2)}(Z) := \sum_{q=0}^n (-1)^q \dim H_{(2)}^{n,q}(Z)$$

*is a bimeromorphic invariant of  $Z$ .*

**Remark.** There are variations of C-G-M conjecture arising from the theory of variations of Hodge structures. In this decade there have been fruitful works appearing on this topic (cf. [44, 42, 22, 8, 20]). As for the statement of the results, see Saper's exposition in this volume.

5. Let us mention other results on the  $L^2$  cohomology groups, which are on pseudoconvex manifolds and therefore more directly related to function theory.

i) Let  $(M, ds^2_M)$  be a Kähler manifold of dimension  $n$  and  $D \subset M$  a relatively compact pseudoconvex domain with  $C^\infty$  smooth boundary. Then the Hodge theory extends to  $D$  as follows.

**Theorem 12** (cf. [38]). *Let  $D$  be as above and let  $q$  be any  $C^\infty$  defining function of  $\partial D$ . If the Levi form of  $q$  along the holomorphic tangent vectors of  $\partial D$  has everywhere at*

least  $n - k$  positive eigenvalues, there exists a complete Kähler metric on  $D$  such that

$$H_{(2)}^{p,q}(D) \cong \begin{cases} H_0^{p,q}(D) & \text{if } p + q \leq n - k \\ H^{p,q}(D) & \text{if } p + q \geq n + k. \end{cases}$$

**Corollary.** Under the above situation,

$$1) \quad H_0^r(D) \cong \bigoplus_{p+q=r} H_0^{p,q}(D) \quad \text{if } r \leq n - k$$

$$H_0^{p,q}(D) \cong \overline{H_0^{q,p}(D)} \quad \text{if } p + q \leq n - k.$$

$$2) \quad H^r(D) \cong \bigoplus_{p+q=r} H^{p,q}(D) \quad \text{if } r \geq n + k$$

$$H^{p,q}(D) \cong \overline{H^{q,p}(D)} \quad \text{if } p + q \geq n + k.$$

- 3) The restriction homomorphisms  $H^{p,q}(D) \rightarrow \lim_{\rightarrow} H^{p,q}(D \setminus K)$  are surjective if  $p + q < n - k$ .

*Remark.* Under somewhat stronger assumptions, results in the above corollary can be proved by different methods. (cf. [28, 29, 26, 12 and 6]).

ii) Donnelly-Fefferman [14] determined the  $L^2$  cohomology of strictly pseudoconvex domains in  $C^n$  with respect to the Bergman metrics.

**Theorem 13.** Let  $D \subset C^n$  be a bounded strictly pseudoconvex domain. Then, with respect to the Bergman metric,

$$\dim H_{(2)}^{p,q}(D) = \begin{cases} 0 & \text{if } p + q \neq n \\ \infty & \text{if } p + q = n \end{cases}$$

We note that the infinite dimensionality of the  $L^2$  cohomology can be proved by an elementary method (cf. [36]), but few things are known about the properties of the  $L^2$  harmonic forms. In case  $p = n$ , the following is known.

**Theorem 14** (cf. [37]). Let  $X$  be a Stein manifold of dimension  $n$ ,  $f$  a bounded holomorphic function on  $X$  such that  $S := \{z \in X; f(z) = 0\}$  has no singular points and  $df|S \neq 0$ , and let  $\varrho_s: H^{n,0}(X \setminus S) \rightarrow H^{n-1,0}(S)$  be the residue homomorphism. Then there exists a bounded linear operator  $I_s: H_{(2)}^{n-1,0}(S) \rightarrow H_{(2)}^{n,0}(X)$  such that  $\varrho_s \cdot (f^{-1}I_s) = \text{id}$ .

**Corollary.** Let  $D \subset C^n$  be a bounded pseudoconvex domain and let  $H \subset C^n$  be a complex hyperplane. Then, every  $L^2$  holomorphic function  $f$  on  $H \cap D$  has an  $L^2$  holomorphic extension to  $D$ , say  $F$  that satisfies  $\|F\| \leq c_D \|f\|$ . Here  $c_D$  is a positive number that depends only on the diameter of  $D$ .

We leave two open questions related to the above results.

1. Under the situation of Theorem 12, let  $x$  and  $y$  be two distinct points in  $D$ . Does there exist an  $L^2$  harmonic  $(p, q)$ -form that separates  $x$  and  $y$  if  $p + q \neq n$ ?

2. Under the situation of Corollary to Theorem 14, let  $g$  be a holomorphic function on  $H \cap D$  which admits a  $C^\infty$  extension to  $D$ , say  $G$ , such that  $\|\bar{\partial}G\| < \infty$ . Then does there exist an  $L^2$  holomorphic extension of  $g$  to  $D$ ?

## References

1. Akizuki, Y., Nakano, S.: Note on Kodaira-Spencer's proof of Lefschetz theorems. Proc. Japan Acad. **30** (1954) 266–272
2. Andreotti, A., Grauert, H.: Théorème de finitude pour la cohomologie des espaces complexes. Bull. Soc. Math. France **90** (1962) 193–259
3. Andreotti, A., Vesentini, E.: Sopra un teorema di Kodaira. Ann. Sci. Norm. Sup. Pisa **15** (1961) 283–309
4. Andreotti, A., Vesentini, E.: Carleman estimates for the Laplace-Beltrami equation on complex manifolds. Publ. Math. IHES **25** (1965) 81–130
5. Arapura, D.: Local cohomology of sheaves of differential forms and Hodge theory. Preprint
6. Bauer, I., Kosarew, S.: On the Hodge spectral sequence for some classes of non-complete algebraic manifolds. Math. Ann. **284** (1989) 577–593
7. Bochner, S.: Curvature and Betti numbers, I, II. Ann. Math. **49** (1948) 379–390; **50** (1949) 77–93
8. Cattani, E., Kaplan, A., Schmidt, W.: Degeneration of Hodge structures. Ann. Math. **123** (1986) 457–535
9. Cheeger, J.: On the Hodge theory of Riemannian pseudomanifolds. Proc. Symp. Pure Math. **36** (1980) 91–146
10. Cheeger, J., Goresky, M., MacPherson, R.:  $L^2$  cohomology and intersection homology for singular algebraic varieties. Seminar on differential geometry. (Ann. Math. Stud. 102.) Princeton Univ. Press 1982, pp. 303–340
11. Demailly, J.P.: Estimations  $L^2$  pour l'opérateur  $\bar{\partial}$  d'un fibré vectoriel holomorphe semi-positif au-dessus d'une variété kähleriennne complète. Ann. Sci. Ec. Norm. Sup. **15** (1982) 457–512
12. Demailly, J.P.: Cohomology of  $q$ -convex spaces in top degrees. Math. Z. **204** (1990) 283–295
13. Diederich, K., Ohsawa, T.: On the parameter dependence of solutions to the  $\bar{\partial}$ -equation. To appear in Math. Ann.
14. Donnelly, H., Fefferman, C.:  $L^2$  cohomology and index theorem for the Bergman metric. Ann. Math. **118** (1984) 593–619
15. Forster, O., Ohsawa, T.: Complete intersections with growth conditions. (Adv. Stud. Pure Math. 10.) Algebraic geometry, Sendai, 1987, pp. 91–104
16. Grauert, H.: On Levi's problem and the imbedding of real-analytic manifolds. Ann. Math. **68** (1958) 460–472
17. Grauert, H., Riemenschneider, O.: Kählersche Mannigfaltigkeiten mit hyper- $q$ -konvexen Rand. Problems in Analysis. Symp. in Honor of S. Bochner. Princeton Univ. Press, 1970, pp. 61–79
18. Hsiang, W.C., Pati, V.:  $L^2$ -cohomology of normal algebraic surfaces. Invent. math. **81** (1985) 395–412
19. Hörmander, L.:  $L^2$  estimates and existence theorems for the  $\bar{\partial}$ -operator. Acta Math. **113** (1965) 89–152
20. Kashiwara, M., Kawai, T.: The Poincaré lemma for variations of polarized Hodge structure. Publ. RIMS, Kyoto Univ. **23** (1987) 345–407
21. Kodaira, K.: On Kähler varieties of restricted type. Ann. Math. **60** (1954) 28–48

22. Looijenga, E.:  $L^2$ -cohomology of locally symmetric varieties. *Comp. Math.* **67** (1988) 3–20
23. Looijenga, E., Rapoport, M.: Weights in the local cohomology of a Baily-Borel compactification. Preprint
24. Nagase, M.: Remarks on the  $L^2$ -cohomology of singular algebraic surfaces. *J. Math. Soc. Japan* **41** (1989) 97–116
25. Nakano, S.: Vanishing theorems for weakly 1-complete manifolds, II. *Publ. RIMS* (1974) 101–110
26. Navarro-Aznar, V.: Sur la théorie de Hodge des variétés algébriques à singularités isolées. *Astérisque* **130** (1985) 272–397
27. Ohsawa, T.: On complete Kähler domains with  $C^1$  boundary. *Publ. RIMS* **10** (1980) 929–940
28. Ohsawa, T.: A reduction theorem for cohomology groups of very strongly  $q$ -convex Kähler manifolds. *Invent. math.* **63** (1981) 335–354
29. Ohsawa, T.: Addendum to: A reduction theorem for cohomology groups of very strongly  $q$ -convex Kähler manifolds. *Invent. math.* **66** (1982) 391–393
30. Ohsawa, T.: Boundary behavior of the Bergman kernel function on pseudoconvex domains. *Publ. RIMS* **20** (1984) 897–902
31. Ohsawa, T.: Vanishing theorems on complete Kähler manifolds. *Publ. RIMS* **20** (1984) 21–38
32. Ohsawa, T.: Hodge spectral sequence on compact Kähler spaces. *Publ. RIMS* **23** (1987) 262–274
33. Ohsawa, T.: Hodge spectral sequence and symmetry on compact Kähler spaces. *Publ. RIMS* **23** (1987) 613–625
34. Ohsawa, T.: Cheeger-Goreski-MacPherson's conjecture for the varieties with isolated singularities. *Math. Z.* **206** (1991) 219–224
35. Ohsawa, T.: Supplement to "Hodge spectral sequence on compact Kähler spaces". To appear in *Publ. RIMS*
36. Ohsawa, T.: On the infinite dimensionality of the middle  $L^2$  cohomology of complex domains. *Publ. RIMS* (1989) 499–502
37. Ohsawa, T., Takegoshi, K.: On the extension of  $L^2$  holomorphic functions. *Math. Z.* **195** (1987) 197–204
38. Ohsawa, T., Takegoshi, K.: Hodge spectral sequence on pseudoconvex domains. *Math. Z.* **197** (1988) 1–12
39. Oka, K.: Domaines finis sans point critique intérieur, *Jap. J. Math.* **27** (1953) 97–155
40. Pardon, Stern, M.:  $L^2 - \bar{\partial}$ -cohomology of complex projective varieties. Preprint
41. Saper, L.:  $L_2$ -cohomology of Kähler varieties with isolated singularities. Preprint
42. Saper, L., Stern, M.:  $L_2$ -cohomology of arithmetic varieties. *Ann. Math.*, to appear
43. Skoda, H.: Morphismes surjectifs de fibrés vectoriels semi-positifs. *Ann. Sci. Ec. Norm. Sup. 4<sup>e</sup> série* **11** (1978) 577–611
44. Zucker, S.:  $L^2$ -cohomology of warped products and arithmetic groups. *Invent. math.* **70** (1982) 169–218



# Differentiability and Measures in Banach Spaces

David Preiss

Department of Mathematics, University College London, London WC1E 6BT, UK

The purpose of this contribution is to give information about new results concerning natural questions about differentiability and measures in real Banach spaces (of infinite but also of finite dimension) and, possibly more importantly, to point out some of the many open problems we are still faced with in this area of research.

## 1. Differentiability

We recall two well known notions.

1. A real valued function  $f$  defined on an open subset  $G$  of a Banach space  $E$  is said to be Fréchet differentiable at a point  $x \in G$  if there is  $f'(x) \in E^*$  such that

$$\lim_{u \rightarrow 0} \frac{|f(x+u) - f(x) - \langle f'(x), u \rangle|}{\|u\|} = 0,$$

$f'(x)$  is called the Fréchet derivative of  $f$  at  $x$ .

2. A real valued function  $f$  defined on an open subset  $G$  of a Banach space  $E$  is said to be Lipschitz on  $G$  if there is a constant  $C$  such that  $|f(x) - f(y)| \leq C\|x - y\|$  whenever  $x, y \in G$ . The smallest such constant  $C$  is denoted by  $\text{Lip}(f)$ .

From the work of Lebesgue (in the one dimensional case) and of Rademacher (in the finite dimensional case) we know that Lipschitz functions on finite dimensional spaces are (Fréchet) differentiable almost everywhere with respect to the Lebesgue measure. Infinite dimensional results of similar nature are known for Gateaux differentiability. (See [1,3,5,6]). These extension are obtained by a linear approximation of the infinite dimensional situation by finite dimensional spaces. However, the question of Fréchet differentiability seems to need a different approach. This might be also seen from many examples of nowhere Fréchet differentiable Lipschitz mappings of a separable Hilbert space into itself, since for such mappings the Gateaux differentiability results mentioned above still hold.

Thus our first result answers a natural question.

**Theorem 1.** *Every Lipschitz function defined on a separable Hilbert space is Fréchet differentiable at least at one point.*

Hilbert spaces are, of course, not the most general spaces in which one would hope for such a result. Indeed, from the extensive investigations of differentiability questions for continuous convex functions (e.g., [13,14]) we know that the result may hold in all Asplund spaces. (A Banach space is said to be an Asplund space if the dual of every its separable subspace is separable.) This generalization of Theorem 1 is given in the following statement.

**Theorem 2.** *Every locally Lipschitz function defined on an open subset of an Asplund space is Fréchet differentiable on a dense subset of its domain.*

The method we use need not be confined to Fréchet differentiability. It applies also to so called  $\mathcal{B}$  derivatives, in the definition of which we require the uniform convergence on the members of a given family  $\mathcal{B}$  of bounded subsets of the Banach space satisfying some mild additional assumptions. (The details can be found in [10].) This gives the most general form of the above differentiability results. (However, a recent Haydon's example of an Asplund space without equivalent smooth norms shows that the deduction of Theorem 2 from Theorem 3 is not straightforward.)

**Theorem 3.** *Let  $E$  be a Banach space admitting an equivalent norm which is  $\mathcal{B}$  differentiable away from the origin. Then every locally Lipschitz function defined on an open subset  $G$  of  $E$  is  $\mathcal{B}$  differentiable on a dense subset of  $G$ .*

These statements, as given, are not satisfactory from the point of view of possible applications. For example, suppose that a Lipschitz function  $f$  on a separable Hilbert space has derivative zero at every point at which it is Fréchet differentiable. We would like to be able to deduce that  $f$  is constant. This can be done, since in all the above results the mean value theorem holds. For example, in case of Theorem 3 we prove that the increment of the function over any segment  $[u, v] \subset G$  is majorized by the supremum of the derivatives in the direction  $v - u$  at points at which the function is  $\mathcal{B}$  differentiable.

The proof of the above results requires new information about Lipschitz functions in finite dimensional spaces. Thus, as a byproduct, we get the following curious statement.

**Theorem 4.** *There is a plane set  $N$  of Lebesgue measure zero such that every Lipschitz function defined on the plane is differentiable at some point of  $N$ .*

To describe a set having such a property is quite easy: Any  $G_\delta$  plane set of Lebesgue measure zero containing all lines passing through two different points with rational coordinates will do. This particular example also suggests the reasons why our proof of Fréchet differentiability results is not straightforward. It combines in some way two notions of smallness of a set: First category (hence the  $G_\delta$  part) and measure zero (hence the lines). It seems to be intuitively clear that a similar mixture is impossible on the line. That this is true has been shown in [2] and [15]: Theorem 4 is false on the line.

### 1.1 Construction of a Point of Fréchet Differentiability

The details of the proof can be found in [10]. Here we just point out the main observations. Because of that we restrict our attention to the proof of Theorem 1 only.

We introduce the directional derivatives of  $f$  by

$$f'(x, e) = \lim_{r \rightarrow 0} \frac{f(x + re) - f(x)}{r},$$

and we denote by  $M$  the set of all pairs  $(x, e) \in E \times E^*$  such that  $\|e\| = 1$  and  $f'(x, e)$  exists.

The first basic observation is that if  $(x, e) \in M$  and  $f'(x, e) = \text{Lip}(f)$  then  $f$  is Fréchet differentiable at  $x$ . Even though such a pair need not exist, this suggests that we might attempt to use a maximizing procedure. Thus our plan is to construct inductively a sequence  $(x_k, e_k) \in M$  so that:

1. The sequence  $x_k$  converges to some  $x$ .
2. The sequence  $e_k$  converges to some  $e$ .
3. The directional derivative  $f'(x, e)$  exists.
4. For the pair  $(x, e)$  some variant of the above observation can be used.

To achieve 1, we simply choose  $x_{k+1}$  close to  $x_k$ . This is based on a local form of our observation, namely, that the equality of  $f'(x, e)$  to the limit of the Lipschitz constants on balls around  $x$  suffices for Fréchet differentiability of  $f$  at  $x$ .

Unfortunately, to get 2 is not so simple. Since requirement 4 forces us to take  $f'(x_k, e_k)$  as large as possible, we cannot at the same time prescribe how close should  $e_k$  be to  $e_{k-1}$ . There is also a different objection we should take into account: If our method worked, we would construct not only a point of differentiability, but also a point at which gradient vector exists. This causes no problem in Hilbert spaces, but is impossible in non-reflexive Asplund spaces. (Every liner functional not attaining its maximum on the unit ball gives an example.) Thus an idea suggests itself: We should change the norm (and the change should depend on  $f$ ), at least in the general case. Recalling that we are constructing a sequence  $e_k$  of unit vectors, and observing that a small change of the norm can drastically change the set of pairs considered for the choice of  $(x_{k+1}, e_{k+1})$ , we find that 2 can be achieved by constructing, together with the sequence  $(x_k, e_k)$ , a sequence of norms  $p_k$ , where  $p_{k+1}$  is the (e.g.,  $l_2$ ) sum of  $p_k$  and of a (small) multiple of the distance to the one dimensional subspace of  $H$  generated by  $e_k$ . Then the conditions  $p_{k+1}(e_{k+1}) = 1$  ( $= p(e_k)$ ) and  $f'(x_{k+1}, e_{k+1}) > f'(x_k, e_k)$  already imply that  $e_{k+1}$  is close to  $e_k$ .

The requirement 3 seems to be the most difficult. To get it, we observe that the problem is essentially one dimensional and requires some method of interchange of limit and derivative. Since we cannot hope to be able to use anything like the uniform convergence of the derivatives, the only possibility seems to be to choose the points at which the increment of the function is approximated by the derivative globally. The following one dimensional lemma says that this can be done.

**Lemma 5.** Suppose that  $a < \xi < b$ ,  $0 < \sigma < 1/4$ , and  $L > 0$  are real numbers,  $h$  is a Lipschitz function defined on  $[a, b]$ ,  $\text{Lip}(h) \leq L$ ,  $h(a) = h(b) = 0$ , and  $h(\xi) \neq 0$ . Then there is a measurable set  $A \subset (a, b)$  such that

1. The Lebesgue measure of the set  $A$  is at least  $\sigma|h(\xi)|/L$ ,
2.  $h'(\tau) \geq \sigma|h(\xi)|/(b - a)$  for every  $\tau \in A$ , and
3.  $|h(t) - h(\tau)| \leq 4(1 + 2\sigma)\sqrt{h'(\tau)L}|t - \tau|$  for every  $\tau \in A$  and every  $t \in [a, b]$ .

The most important third statement of the lemma says that the approximation of the increment of the function by its derivative at the point  $\tau$  is “globally good” in the whole interval  $[a, b]$ . The second statement just says that the derivative at  $\tau$  increased as much as we could hope for. From the first statement we just use that  $\tau$  can be chosen sufficiently far from the end points. This is needed in order to get a bilateral approximation.

Because in the first statement of Lemma 5 we do not have to speak about measure, the Lemma can be formulated without the notion of the Lebesgue measure. We can then try to prove it without any use of measure theory. This sounds difficult, since we also claim that  $h$  is differentiable at  $\tau$ . But we can also replace the derivatives by lower derivatives and get a version of the lemma that really can be proved without any use of measure theory. Surprisingly enough, this statement then easily implies that Lipschitz functions on the real line have at least one point of differentiability. Though I did not follow this way, since to use the Lebesgue measure and maximal operator technique turned out to be much easier, these remarks suggest that the proof of differentiability discussed here is different from the usual measure theoretic proofs.

Having done this, we can already imagine how to construct the sequence  $(x_k, e_k)$  so that 3 holds: We will choose  $(x_{k+1}, e_{k+1})$  so that the approximation of the the increment of the function by its directional derivative at the point  $x_{k+1}$  in the direction  $e_{k+1}$  is “globally good” on the whole line through  $x_{k+1}$  in the direction  $e_{k+1}$ .

However, the previous choice implies that our construction will lead to a pair  $(x, e)$  for which the equality  $f'(x, e) = \text{Lip}(f)$  is quite far from being true. Hence to achieve 4 we need to improve upon our main observation. We first reformulate this observation as:

A Lipschitz function  $f$  on a Hilbert space is Fréchet differentiable at  $x$  if there is a unit vector  $e$  such that  $f'(x, e)$  exists and

$$\limsup_{\delta \searrow 0} \{f'(\tilde{x}, \tilde{e}); (\tilde{x}, \tilde{e}) \in M \text{ and } \|\tilde{x} - x\| < \delta\} \leq f'(x, e).$$

A simple proof of this statement together with Lemma 5 gives the following differentiability criterion, which we formulate in the most general situation.

**Theorem 6.** Suppose that  $E$  is a Banach space,  $x_0 \in E$ ,  $e_0 \in E$ ,  $\|e_0\| = 1$ , and that  $f$  is a Lipschitz function defined on  $E$  such that  $f'(x_0, e_0)$  exists. Let  $M$  denote the set of all pairs  $(x, e) \in E \times \{e \in E; \|e\| = 1\}$  such that  $f'(x, e)$  exists,  $f'(x, e) \geq f'(x_0, e_0)$ , and

$$|(f(x + te_0) - f(x)) - (f(x_0 + te_0) - f(x_0))| \leq 6|t|\sqrt{(f'(x, e) - f'(x_0, e_0))\text{Lip}(f)}$$

for every  $t \in \mathbb{R}$ .

Then, if the norm is  $\mathcal{B}$  differentiable at  $e_0$ , and if

$$\lim_{\delta \searrow 0} \sup\{f'(x, e); (x, e) \in M \text{ and } \|x - x_0\| \leq \delta\} \leq f'(x_0, e_0),$$

$f$  is  $\mathcal{B}$  differentiable at  $x_0$ .

Now, the way of constructing the sequences  $(x_k, e_k)$  and  $p_k$  is more or less clear. We always pick up the next pair from the set  $M$  described in the previous Theorem. The additional requirement is only that  $f'(x_{k+1}, e_{k+1})$  is very close to the supremum of the directional derivatives  $f'(x, e)$  for  $(x, e) \in M$ . Then we define the norm  $p_{k+1}$  and continue our construction. Though we still have to be quite careful and make some technical estimates, since, for example, the set  $M$  from the previous Theorem depends upon the choice of the norm, this construction leads to a sequence satisfying all our requirements.

## 1.2 Problems

From the previous discussion it is clear that the theory of differentiability still abounds with open problems. I would just like to point out the following two.

**Problem 7.** Does every pair of Lipschitz functions on a separable Hilbert space have a common point of differentiability?

**Problem 8.** For which finite Borel measures in separable Banach spaces is it true that every Lipschitz function is differentiable almost everywhere?

The second problem is purely finite dimensional since such measures do not exist in infinite dimensional spaces. (See [12].) The answer is not known in the plane (or in any higher dimensional space). In the one dimensional case the required measures are precisely those that are absolutely continuous with respect to the Lebesgue measure. In spite of Theorem 4 I do not know any example that would show that this is not true in all finite dimensional spaces.

## 2. Measures

The question whether measures on separable Banach spaces are determined by their values on balls has been around since R. O. Davies [4] published his beautiful example of two different probability measures on a compact metric space that agree on all balls. Together with J. Tišer [11] we recently answered it by proving:

**Theorem 9.** Whenever two finite Borel measures in a separable Banach space agree on all balls, then they agree.

To prove this statement, we first use the Fourier transform to reduce the problem to showing that the measures agree on all halfspaces. Then, by blowing up balls, we come to the situation when the halfspace contains a nonempty open cone  $C$  on every translate of which the measures agree. An approximation argument (or a differentiability result from [7]) reduces the problem further to the case when  $\overline{C} \cap -\overline{C}$  is a subspace of finite codimension. Thus we can pass to the factor space and we have to solve the corresponding problem in finite dimensional spaces: Do we know the measure of a halfspace provided we know the measure of each translate of a nonempty open cone contained in it? Since this turned out to be true, our approach has been successful.

Instead of giving further details, it might be more interesting to point out some examples. The motivation for them comes from the Besicovitch-Morse differentiability theorem, which is a much stronger statement than that measures in finite dimensional normed spaces are determined by their values on balls:

For every (locally) finite Borel measure  $\mu$  in a finite dimensional Banach space and for every  $\mu$  integrable function  $f$  the limit

$$\lim_{r \searrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} f(u) d\mu(u) \quad (1)$$

exists and equals  $f(x)$  for  $\mu$  almost every  $x$ .

As a corollary of this statement one can prove that, if  $\mu$  and  $\nu$  are two finite Borel measures in a finite dimensional Banach space satisfying  $\mu(B) \geq \nu(B)$  for every ball  $B$  then  $\mu \geq \nu$ .

*Example 1* ([9]). There is a Gaussian measure  $\gamma_1$  in  $l_2$  and a  $\gamma_1$  integrable function  $f$  such that the limit in (1) is infinite uniformly for  $x \in l_2$ , i.e.,

$$\liminf_{r \searrow 0} \frac{1}{\gamma_1(B(x, r))} \int_{B(x, r)} f(u) d\gamma_1(u) = \infty.$$

*Example 2* ([8]). There is a Gaussian measure  $\gamma_2$  in  $l_2$  and a bounded  $\gamma_2$  measurable function  $f$  such that, for  $\mu$  almost every  $x$ , the limit in (1) does not exist.

*Example 3* (J. Tišer). There is a non-degenerated Gaussian measure  $\gamma_3$  in  $l_2$  such that (1) holds for every  $f \in L_p(\gamma_3)$ ,  $p > 1$ .

*Example 4.* In a separable Hilbert space the statement " $\mu(B) \geq \nu(B)$  for balls with radius less than one implies  $\mu \geq \nu$ " holds if and only if the dimension of the space is finite.

*Example 5.* In a separable Hilbert space the statement " $\mu(B) \geq \nu(B)$  for balls with radius greater than one implies  $\mu \geq \nu$ " holds if and only if the dimension of the space is infinite.

*Example 6.* In the  $l_\infty$  sum of a separable Hilbert space with the line there are measures  $\mu$  and  $\nu$  such that  $\mu \not\geq \nu$  but  $\mu(B) \geq \nu(B)$  for all balls.

However, in spite of the above result and examples, the investigation of the behaviour of measures on balls cannot be considered as finished. For example, the following question is still far from being answered.

**Problem 10.** Are finite Borel measures in separable Banach spaces determined by their values on balls with radii less than one?

## References

1. Aronszajn, N.: Differentiability of Lipschitz functions in Banach spaces. *Studia Math.* **57** (1976) 147–160
2. Choquet, G.: Applications des propriétés descriptives de la fonction contingent à la théorie des fonctions de variable réelle et à la géométrie différentielle des variétés cartésiennes. *J. Math. Pures Appl.* **26** (1947) 115–226
3. Christensen, J. P. R.: Measure theoretic zero sets in infinitely dimensional spaces and application to differentiability of Lipschitz mappings. *Actes du Deuxieme Colloque d'Analyse Fonctionnelle de Bordeaux* **2** (1973) 29–39
4. Davies, R. O.: Measures not approximable or not specifiable by means of balls. *Mathematika* **18** (1971) 157–160
5. Mankiewicz, P.: On the differentiability of Lipschitz mappings in Fréchet spaces. *Studia Math.* **45** (1973) 15–29
6. Phelps, R. R.: Gaussian null sets and differentiability of Lipschitz maps on Banach spaces. *Pacific J. Math.* **77** (1978) 523–531
7. Preiss, D.: Almost differentiability of convex functions on Banach spaces and determination of measures by their values on balls. In: *Proc. Conf. Geometry of Banach Spaces (Strobl 1989)* (to appear)
8. Preiss, D.: Gaussian measures and the density theorem. *Comment. Math. Univ. Carolinae* **22** (1981) 181–193
9. Preiss, D.: Differentiation of measures in infinitely dimensional spaces. In: *Proc. Conf. Topology and Measure III*, pages 201–207, Greifswald, 1982
10. Preiss, D.: Differentiability of Lipschitz functions on Banach spaces. *J. Functional Anal.* **91** (1990) 312–345
11. Preiss, D., Tišer, J.: Measures on Banach spaces are determined by their values on balls (to appear)
12. Preiss, D., Zajíček, L.: Fréchet differentiation of convex functions in a Banach space with a separable dual. *Proc. Amer. Math. Soc.* **91** (1984) 202–204
13. Stegall, Ch.: The duality between Asplund spaces and spaces with the Radon-Nikodym property. *Israel J. Math.* **29** (1978) 408–412
14. Stegall, Ch.: The Radon-Nikodym property in conjugate Banach spaces II. *Trans. Amer. Math. Soc.* **264** (1981) 507–519
15. Zahorski, Z.: Sur l'ensemble des points de non-derivabilité d'une fonction continue. *Bull. Soc. Math. France* **74** (1946) 147–178



# The Limit Element in the Configuration Algebra for a Discrete Group: A précis

Kyoji Saito

Research Institute for Mathematical Sciences, Kyoto University, Kyoto 606, Japan

## §1. Introduction

(1.1) We construct certain infinitely generated Hopf algebra, which we call the *configuration algebra*. The algebra is generated by isomorphism classes of colored oriented finite graphs and is completed with respect to an adic topology filtered by the cardinality of graphs. Therefore certain limit process is admitted inside the algebra, which is absolutely necessary for the application explained in (1.2).

(1.2) The original attempt of the work, which is under investigation, is to apply the algebra for a construction of certain modular function on the moduli of discrete groups [8]. Let us explain this. Consider a finitely generated group  $\Gamma$ . By fixing a generator system,  $\Gamma$  gets naturally colored oriented graph structure, the Caylay graph. Let  $\Gamma_n$  be the subgraph of  $\Gamma$  consisting of elements of length  $\leq n$  and let  $\mathcal{A}(\Gamma_n)$  be the formal sum in the algebra of all non-void subgraphs of  $\Gamma_n$ . Then the following limit process is justified in the configuration algebra:

$$\omega_\Gamma := \lim_{n \rightarrow \infty} \log \left( (1 + \mathcal{A}(\Gamma_n))^{1/\#\Gamma_n} \right).$$

Namely, we show that these elements  $\log(1 + \mathcal{A}(\Gamma_n))$  becomes Lielike in the algebra. By a use of certain basis  $\{\varphi(S)\}_S$  (where  $S$  runs the set of isomorphism classes of connected graphs) for the space  $\mathcal{L}_{\mathbb{R}}$  of Lielike elements (§8), the above limit element  $\omega_\Gamma$  is developed as follows.

$$\omega_\Gamma = \sum_{S \in \text{Conf}_0} \varphi(S) \cdot \lim_{t \rightarrow r} \frac{PM(S, t)}{PM(pt, t)}.$$

Here  $PM(S, t) := \sum_{n=0}^{\infty} A(S, \Gamma_n) t^n$  is the generating function for  $A(S, \Gamma_n) := \#\{\text{subgraphs of } \Gamma_n \text{ isomorphic to } S\}$ ,  $pt$  = the graph of one vertex, and  $r > 0$  is the radius of convergence of  $PM(pt, t)$ . For a wide class of groups  $\Gamma$ , including hyperbolic groups and certain automatic groups (cf. [4, 2]), the ratio  $PM(S, t)/PM(pt, t)$  extends to a rational function in  $t$  which is *regular* at  $t = r$  (see §10). So one obtains a final formula,

$$\omega_r = \sum_{S \in \text{Conf}_0} \varphi(S) \cdot \frac{PM(S, t)}{PM(pt, t)} \Big|_{t=r}$$

where  $r := 1 / \limsup_{n \rightarrow \infty} \sqrt{\#\Gamma_n}$  is the first place of pole of  $PM(pt, t)$  on the real axis and hence is a real algebraic number and  $\omega_r \in \mathcal{L}_{\mathbb{Q}(r)}$ .

To obtain the modular function, we need further study of representations of the Configuration algebra, which is a subject of another paper.

The §§2–4 treat generality on graphs and construct the configuration algebra. The §§5–9 treat Lielike elements and grouplike elements of the algebra. The §10 treats the limit elements in the algebra.

## §2. Graphs and Covering Coefficients

Some basic combinatorial rules for graphs are discussed in this paragraph.

**(2.1) Definition.** 1. A pair  $(\Gamma, B)$  is called a *graph*, if  $\Gamma$  is a set of vertices and  $B$  (called the set of edges) is a subset of  $\Gamma \times \Gamma \setminus \Delta$  with  $\sigma(B) = B$ , where  $\sigma$  is the involution  $\sigma(\alpha, \beta) = (\beta, \alpha)$  and  $\Delta$  is the diagonal. A graph is *connected*, if it is connected as a simplicial complex. A graph is *finite*, if  $\#\Gamma < \infty$ .

2. An *isomorphism* of graphs is a bijection of vertices inducing a bijection of edges. Any subset  $\mathbb{S}$  of  $\Gamma$  is a *subgraph* by taking  $B \cap (\mathbb{S} \times \mathbb{S})$  as edges for  $\mathbb{S}$ . The word “subgraph” is used only in this sense.

3. A graph is called *colored oriented*, if there exists a finite set  $G$ , called a coloring set, with an involution  $\sigma_G : G \rightarrow G$  and a map  $c : B \rightarrow G$  which is equivariant with the involutions:  $c \circ \sigma = \sigma_G \circ c$ . Isomorphisms and subgraphs of colored oriented graphs are defined as compatible with  $c$ .

*Example.* A group  $\Gamma$  with a finite generator system  $G$  with  $G = G^{-1}$  and  $e \notin G$  carries a colored oriented graph structure:  $B := \{(\gamma, \delta) \in \Gamma^2 : \gamma^{-1}\delta \in G\}$ ,  $c(\gamma, \delta) = \gamma^{-1}\delta$  and  $\sigma_G(g) = g^{-1}$ . The graph is called a *Cayley graph*.

**(2.2)** In all what follow, we fix an increasing sequence  $\{G_p\}_{p \in \mathbb{N}}$  of coloring set. Associated to that the set of configurations is defined by

$$\begin{aligned} \text{Conf}^{pq} := \{&\mathbb{S} : \mathbb{S} \text{ is a finite colored oriented graph for } G_p \text{ such that} \\ &\text{number of edges at a vertex is at most } q\}/\text{isomorphism}. \end{aligned}$$

$$\text{Conf} := \cup_{p,q} \text{Conf}^{pq},$$

$$\text{Conf}_0 := \{S \in \text{Conf} : S \text{ is connected}\}, \quad \text{Conf}_+ := \text{Conf} \setminus \{[\phi]\}.$$

An isomorphisms class of a graph  $\mathbb{S}$  is denoted by  $[\mathbb{S}]$ . The set of vertices of  $\mathbb{S}$  is denoted by  $|\mathbb{S}|$ .

**(2.3)** The Conf has an abelian semigroup structure with a partial ordering:

$$[\mathbb{S}] \cdot [\mathbb{T}] := [\mathbb{S} \amalg \mathbb{T}] \quad \text{for } [\mathbb{S}], [\mathbb{T}] \in \text{Conf},$$

$S \leq T \stackrel{\text{def}}{\iff} \text{There are graphs } \mathbf{S} \text{ and } \mathbf{T} \text{ with } S = [\mathbf{S}], T = [\mathbf{T}] \text{ and } \mathbf{S} \subset \mathbf{T}.$   
 Here  $1 = [\phi]$  is the minimal element in  $\text{Conf}$ . One has  $\text{Conf} \simeq \text{Conf}_0^{\mathbb{Z}_{\geq 0}}$ .

(2.4) For  $S_1, \dots, S_m$  and  $S \in \text{Conf}$ , we define a numerical invariant, which we call the *covering coefficient*:

$$\binom{S_1, \dots, S_m}{S} := \#\{(\mathbf{S}_1, \dots, \mathbf{S}_m) : \mathbf{S}_i \subset \mathbf{S} \text{ s.t. } [\mathbf{S}_i] = S_i \ (i = 1, \dots, m) \text{ and } \bigcup_{i=1}^m |\mathbf{S}_i| = |\mathbf{S}|$$

Here  $\mathbf{S}$  is a graph s.t.  $[\mathbf{S}] = S$ . The definition does not depend on  $\mathbf{S}$ .

(2.5) We list some elementary properties of covering coefficients.

- i)  $\binom{S_1, \dots, S_m}{S} = 0$  unless  $S_i \leq S$  for  $i = 1, \dots, m$  and  $\sum \#S_i \geq \#S$ .
- ii)  $\binom{S_1, \dots, S_m}{S}$  is invariant by the permutation of  $S_i$ 's.
- iii) For  $1 \leq i \leq m$ , one has an elimination rule:

$$\binom{S_1, \dots, S_{i-1}, \phi, S_{i+1}, \dots, S_m}{S} = \binom{S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_m}{S}.$$

iv) For the case  $m = 1$ ,  $\binom{T}{S} = \begin{cases} 1 & \text{if } S = T, \\ 0 & \text{else.} \end{cases}$

v) For the case  $S = \phi$ ,  $\binom{S_1, \dots, S_m}{\phi} = \begin{cases} 1 & \text{if } \cup S_i = \phi, \\ 0 & \text{else.} \end{cases}$

(2.6) **Composition Rule.** For  $S_1, \dots, S_m, T_1, \dots, T_n$  and  $S \in \text{Conf}$ , one has

$$\sum_{U \in \text{Conf}} \binom{S_1, \dots, S_m}{U} \binom{U, T_1, \dots, T_n}{S} = \binom{S_1, \dots, S_m, T_1, \dots, T_n}{S}.$$

(2.7) **Decomposition Rule.** For  $S_1, \dots, S_m, U$  and  $V \in \text{Conf}$ , one has

$$\binom{S_1, \dots, S_m}{U \cdot V} = \sum_{S_1 = R_1 \cdot T_1} \dots \sum_{S_m = R_m \cdot T_m} \binom{R_1, \dots, R_m}{U} \binom{T_1, \dots, T_m}{V}.$$

Here  $R_i$  and  $T_i$  run over  $\text{Conf}$  for all possible decompositions of  $S_i$  ( $i = 1, \dots, m$ ).

### §3. Configuration Algebra

(3.1) *Algebra Structure.* The free abelian group  $\mathbb{Z} \cdot \text{Conf}$  generated by  $\text{Conf}$  carries an algebra structure by a use of the semigroup structure on  $\text{Conf}$  (2.3) as the product (recall that  $[\phi] = 1$ ). It is isomorphic to the free polynomial algebra generated by  $\text{Conf}_0$ . The algebra is graded by taking  $\deg(S) := \#(S)$  for  $S \in \text{Conf}$ , since the aditivity:  $\#(S \cdot T) = \#(S) + \#(T)$ .

(3.2) *Adic Topology.* For  $n \geq 0$ , let us define an ideal in  $\mathbb{Z} \cdot \text{Conf}^{pq}$

$$\mathcal{J}_n^{pq} := \text{the ideal generated by } \{S \in \text{Conf}^{pq} : \#(S) \geq n\}.$$

Taking  $\mathcal{J}_n^{pq}$  as a fundamental system of neighbourhoods of 0, we define the adic topology. The completion w.r.t. the adic topology is denoted by

$$\mathbb{Z}[[\text{Conf}^{pq}]] := \varprojlim_n (\mathbb{Z} \cdot \text{Conf}^{pq}) / \mathcal{J}_n^{pq}.$$

We put  $\mathbb{Z}[[\text{Conf}]] := \varinjlim \mathbb{Z}[[\text{Conf}^{pq}]]$ , even it is a confusing notation. For a commutative algebra  $\mathbb{A}$  with a unit, we put  $\mathbb{A}[[\text{Conf}]] := \mathbb{A} \otimes_{\mathbb{Z}} \mathbb{Z}[[\text{Conf}]]$  and call it the *configuration algebra* with coefficient in  $\mathbb{A}$ . Since  $\mathbb{A}[[\text{Conf}^{pq}]] \simeq \prod_{S \in \text{Conf}^{pq}} \mathbb{A} \cdot S$ , any element  $f$  of the algebra is expressed by a formal sum :  $f = \sum_{S \in \text{Conf}^{pq}} S \cdot f_S$  for a suitable  $p, q > 0$  and constants  $f_S \in \mathbb{A}$ . Put,  $\text{Supp}(f) := \{S \in \text{Conf} : f_S \neq 0\}$ . A series  $f$  is said to be *of finite type*, if  $\text{Supp}(f)$  is contained in a finitely generated subsemigroup of  $\text{Conf}$ .

(3.3) A subset  $P$  of  $\text{Conf}$  is called *saturated*, if its *saturation*  $\tilde{P} := \{T \in \text{Conf} : \exists S \in P \text{ s.t. } T \leq S\}$  coincides itself. A subalgebra generated by a saturated subset is called a *saturated subalgebra* of the configuration algebra. Let  $\Gamma$  be a Caylay graph of an infinite group and let  $P_{\Gamma}$  be the set of isomorphism classes of all finite subgraphs of  $\Gamma$ . Then  $P_{\Gamma}$  is a saturated semigroup contained in  $\text{Conf}^{pq}$  for some  $p, q \geq 0$ .  $\mathbb{A}[[P_{\Gamma}]]$  is the configuration algebra for  $(\Gamma, G)$ .

(3.4) *Exponential and Logarithmic Maps.* Assume  $\mathbb{Q} \subset \mathbb{A}$ . We define maps:

$$\exp(\mathcal{M}) := \sum_{n=0}^{\infty} \frac{1}{n!} \mathcal{M}^n \quad \text{and} \quad \log(1 + \mathcal{A}) := \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} \mathcal{A}^n$$

for  $\mathcal{M}$  and  $\mathcal{A} \in \mathbb{A}[[\text{Conf}]]_+$  := the augmentation ideal generated by  $\text{Conf}_+$ .

## §4. The Hopf Structure on the Configuration Algebra

(4.1) *Coprodutct  $\Phi_m$ .* For an  $m \in \mathbb{N}$  and  $U \in \text{Conf}$ , define a map

$$\Phi_m(U) := \sum_{S_1 \in \text{Conf}} \cdots \sum_{S_m \in \text{Conf}} \binom{S_1, \dots, S_m}{U} S_1 \otimes \dots \otimes S_m.$$

Owing to the decomposition rule (2.7) and (2.5) i), one has

$$\begin{aligned} \Phi_m(U \cdot V) &= \Phi_m(U) \cdot \Phi_m(V) && \text{for } U, V \in \text{Conf}, \\ \Phi_m(\mathcal{J}_n^{pq}) &\subset \sum_{n_1 + \dots + n_m \geq n} \mathcal{J}_{n_1}^{pq} \otimes \dots \otimes \mathcal{J}_{n_m}^{pq} && \text{for } p, q, n, m \in \mathbb{Z}_{\geq 0}. \end{aligned}$$

Thus  $\Phi_m$  extends to a homomorphism from the configuration algebra to its completed  $m$ -tensor product, denoted by  $\Phi_m$  again, which we call the *coproduct*. The symmetric group  $\mathfrak{S}_m$  acts on the range of  $\Phi_m$  by permutation. The image

of  $\Phi_m$  is  $\mathfrak{S}_m$ -invariant, because of (2.5) ii). We shall call this property the *cocommutativity* of  $\Phi_m$ .

The composition rule (2.6) implies, for  $m, n \in \mathbb{N}$ ,

$$(1 \otimes \cdots \otimes 1 \otimes \Phi_m) \circ \Phi_{n+1} = \Phi_{m+n}.$$

Combining this fact with the cocommutativity, we see that any composites of  $\Phi$  is expressed again by some  $\Phi_m$ . We call this the *co-associativity* of  $\Phi_m$ .

**(4.2)** The *augmentation map* for the algebra is defined as  $\mathbb{Z}[[\text{Conf}]] \rightarrow \mathbb{Z}$ , ( $S \in \text{Conf}_+ \mapsto 0$ ,  $[\phi] \mapsto 1$ ). It is a *co-unit* in the sense:  $(\text{aug} \cdot \text{id}) \circ \Phi_2 = \text{id}$ .

**(4.3)** *The involution map*  $\iota$ . The following Lemma is nontrivial.

**Lemma (Existense of Co-inverse).** *There exists a unique algebra automorphism*

$$\iota : \mathbb{Z}[[\text{Conf}]] \rightarrow \mathbb{Z}[[\text{Conf}]]$$

such that

- i)  $\iota$  is involutive. That is:  $\iota^2 = \text{id}$ .
- ii)  $\iota$  is the co-inverse map w.r.t. the coproduct  $\Phi_2$ . That is:  $(\iota \cdot \text{id}) \circ \Phi_2 = \text{aug}$ .
- iii)  $\iota$  leaves any saturated subalgebras of  $\mathbb{Z}[[\text{Conf}]]$  invariant.
- iv)  $\iota$  is continuous w.r.t. the adic topology. In particular,  $\text{aug} \circ \iota = \text{aug}$ .

*Remark.* There is other coalgebra structure studied in combinatorics (Rota [5]).

## §5. The Growth Functions for Configurations

In this and next paragraphs, we introduce some basic elements  $1 + \mathcal{A}(T)$  and  $\mathcal{M}(T)$  ( $T \in \text{Conf}$ ) of the configuration algebra, which turn out to be grouplike or Lie-like respectively.

**(5.1)** *Growth Functions.* Let  $S$  and  $T \in \text{Conf}$  be given. Fix a graph  $\mathbb{T}$  with  $[\mathbb{T}] = T$  and put

$$\begin{aligned} \mathbf{A}(S, \mathbb{T}) &:= \#\{\mathbf{S} : \mathbf{S} \subset \mathbb{T} \text{ such that } [\mathbf{S}] = S\}, \\ A(S, T) &:= \# \mathbf{A}(S, \mathbb{T}). \end{aligned}$$

We call  $A(S, T)$  the growth function. By definition  $A(S, T_1 \cdot T_2) = A(S, T_1) + A(S, T_2)$  for  $S \in \text{Conf}_0$  and  $T_i \in \text{Conf}$ .

Let us introduce an element of  $\mathbb{Z} \cdot \text{Conf}$ :

$$1 + \mathcal{A}(T) := 1 + \sum_{S \in \text{Conf}_+} S \cdot A(S, T) = \sum_{\mathbf{S} \in 2^{\mathbb{T}}} [\mathbf{S}],$$

where  $2^{\mathbb{T}} := \{\text{subgraphs of } \mathbb{T}\}$ . Obvious from definition, for  $T_1$  and  $T_2 \in \text{Conf}$

$$(1 + \mathcal{A}(T_1 \cdot T_2)) = (1 + \mathcal{A}(T_1))(1 + \mathcal{A}(T_2)).$$

(5.2) *A numerical approximation* of the growth function.

**Lemma.** *For  $S \in \text{Conf}$  and for  $T \in \text{Conf}^{pq}$ , one has*

$$\mathcal{A}(S, T) \leq \frac{1}{\# \text{Aut}(S)} \cdot (\#T)^{n(S)} \cdot q^{\#S - n(S)}.$$

(Here  $n(S) := \#$  of connected component of  $S$ .)

(5.3) *Product expansion formula* for growth functions. The following basic formula can be shown by a combinatorial method.

**Lemma.** *Let  $S_1, \dots, S_m$  ( $m \geq 1$ ) and  $T \in \text{Conf}$  be given. Then,*

$$\prod_{i=1}^m A(S_i, T) = \sum_{S \in \text{Conf}} \binom{S_1, \dots, S_m}{S} A(S, T).$$

(5.4) *Grouplike property* of the growth function.

An element  $g \in \mathbb{A}[[\text{Conf}]] \setminus \{0\}$  is called grouplike [8], if it satisfies  $\Phi_m(g) = g \otimes \cdots \otimes g$  for  $\forall m \in \mathbb{N}$ . We put

$$\begin{aligned} G_{\mathbb{A}} &:= \{\text{all grouplike elements in } \mathbb{A}[[\text{Conf}]]\}, \\ G_{\mathbb{A}, \text{finite}} &:= \{g \in G_{\mathbb{A}} : g \text{ is of finite type}\}. \end{aligned}$$

The (5.3) Lemma implies the following Lemma.

**Lemma (Grouplike Property).** *For all  $T \in \text{Conf}$ , one has  $1 + \mathcal{A}(T) \in G_{\mathbb{Z}, \text{finite}}$ . That is: for any  $m \in \mathbb{N}$  and  $T \in \text{Conf}$ , one has*

$$(1 + \mathcal{A}(T)) \otimes \cdots \otimes (1 + \mathcal{A}(T)) = \Phi_m(1 + \mathcal{A}(T)).$$

**Corollary.** 
$$(1 + \iota(\mathcal{A}(T))) (1 + \mathcal{A}(T)) = 1 \quad \text{for } T \in \text{Conf},$$
  

$$\Phi_m \circ \iota = (\iota \otimes \cdots \otimes \iota) \circ \Phi_m \quad \text{for } m \in \mathbb{N}.$$

## §6. The Logarithmic Growth Function

(6.1) For  $S \in \text{Conf}$ , define the logarithmic growth function by:

$$\mathcal{M}(T) := \log(1 + \mathcal{A}(T))$$

(cf. (3.4). Develop  $\mathcal{M}(T)$  in a series

$$\mathcal{M}(T) = \sum_{S \in \text{Conf}} S \cdot M(S, T).$$

By definition  $M(\phi, T) := 0$  for  $T \in \text{Conf}$  and  $A(S, T) = M(S, T)$  for  $S \in \text{Conf}_0$ . The multiplicativity of  $\mathcal{A}(T)$  (5.1) implies the additivity:

$$\mathcal{M}(T_1 \cdot T_2) = \mathcal{M}(T_1) + \mathcal{M}(T_2)$$

for  $T_i \in \text{Conf}$  and hence the additivity:

$$M(S, T_1 \cdot T_2) = M(S, T_1) + M(S, T_2)$$

for  $S$  and  $T_i \in \text{Conf}$ .

### (6.2) The linear dependence relations on logarithmic functions.

The polynomial relation (5.4) Lemma implies a linear relation:

**Lemma (Lielike Property).** For  $T \in \text{Conf}$  and  $m \geq 1$ ,

$$\sum_{i=1}^m 1 \otimes \cdots \otimes 1 \overset{i\text{-th}}{\mathcal{M}}(T) \otimes 1 \otimes \cdots \otimes 1 = \Phi_m(\mathcal{M}(T)).$$

**Corollary.** Let  $S_1, \dots, S_m \in \text{Conf}_+$  for  $m \geq 2$  and  $T \in \text{Conf}$  be given. Then,

$$\sum_{S \in \text{Conf}} \binom{S_1, \dots, S_m}{S} M(S, T) = 0.$$

*Remark.* The linear dependence relations among  $M(S, T)$ 's for  $S \in \text{Conf}$  is the key fact in all this paper. The Hopf algebra structure is introduced to deduce the relation. We shall solve this linear relation in (6.2) by a use of kabi coefficients introduced in §7.

(6.3) An element  $\mathcal{M}$  of the configuration algebra satisfying the relation (6.2) is called *Lielike* ([7]). Let us fix notations:

$$\mathcal{L}_{\mathbb{A}}^{pq} := \{\text{all Lie-like elements in } \mathbb{A}[[\text{Conf}^{pq}]]\},$$

and

$$\mathcal{L}_{\mathbb{A}, \text{finite}}^{pq} := \{M \in \mathcal{L}_{\mathbb{A}}^{pq} : M \text{ is of finite type}\}.$$

We denote also:  $\mathcal{L}_{\mathbb{A}} := \cup_{p,q} \mathcal{L}_{\mathbb{A}}^{pq}$  and  $\mathcal{L}_{\mathbb{A}, \text{finite}} := \cup_{p,q} \mathcal{L}_{\mathbb{A}, \text{finite}}^{pq}$ .

## §7. Kabi Coefficients

(7.1) **Definition.** 1. A pair  $(\mathbb{S}, \mathbb{U})$  of a graph  $\mathbb{U}$  and its subgraph  $\mathbb{S}$  is called a *kabi* over  $\mathbb{S}$ , if any vertex of  $\mathbb{U} \setminus \mathbb{S}$  is connected to a vertex of  $\mathbb{S}$  through an edge.

2. Let  $U \in \text{Conf}_0$  and let  $\mathbb{U}$  be a graph with  $[\mathbb{U}] \subseteq U$ . For  $S \in \text{Conf}_0$ , put

$$\mathbb{K}(S, \mathbb{U}) := \{\mathbb{S} : \mathbb{S} \subset \mathbb{U} \text{ s.t. } [\mathbb{S}] \subseteq S \text{ and } (\mathbb{S}, \mathbb{U}) \text{ is a kabi}\},$$

$$K(S, U) := \#\mathbb{K}(S, \mathbb{U}).$$

We call  $K(S, U)$  a *Kabi-coefficient*. Its definition does not depend on a choice of  $\mathbb{U}$ . From definition,  $K(S, U) = 0$  for  $S \not\leq U$ , and  $K(S, S) = 1$  for  $S \in \text{Conf}_0$ . The word “kabi” means “mold” in Japanese.

(7.2) **Lemma (The Inversion Formula of Kabi).** *For  $S, T \in \text{Conf}_0$ , one has*

$$(1) \quad \sum_{U \in \text{Conf}_0} (-1)^{\#U - \#S} K(S, U) \cdot A(U, T) = \delta(S, T).$$

$$(2) \quad \sum_{U \in \text{Conf}_0} (-1)^{\#U - \#T} A(S, U) \cdot K(U, T) = \delta(S, T).$$

Specializing (2), one gets:  $\sum_{U \in \text{Conf}_0} (-1)^{\#U} \#U \cdot K(U, T) = -\delta(pt, T)$ .

(7.3) We remark the boundedness of non-zero entries of  $K$ .

*Assertion.* Let  $T \in \text{Conf}_0^{pq}$ . Then  $K(S, T) = 0$ , unless  $\#T \leq \#S \cdot (q-1) + 2$ .

## §8. Lielike Elements $\mathcal{L}_{\mathbb{A}}$

(8.1) *Bases of  $\mathcal{L}_{\mathbb{A}, \text{finite}}$  and  $\mathcal{L}_{\mathbb{A}}$ .*

**Lemma.** Let  $\mathbb{A}$  be an algebra containing  $\mathbb{Q}$ . Then,

i) The system  $\{\mathcal{M}(T)\}_{T \in \text{Conf}_0}$  give a  $\mathbb{A}$ -free basis for  $\mathcal{L}_{\mathbb{A}, \text{finite}}$ .

$$\mathcal{L}_{\mathbb{A}, \text{finite}} \simeq \bigoplus_{S \in \text{Conf}_0} \mathbb{A} \cdot \mathcal{M}(S).$$

We introduce another system  $\{\varphi(S) \in \mathcal{L}_{\mathbb{A}, \text{finite}} \cap \mathcal{J}_{\#S}\}_{S \in \text{Conf}_0}$  of  $\mathbb{A}$  basis of  $\mathcal{L}_{\mathbb{A}, \text{finite}}$  (recall the inversion formula (7.2)).

$$\begin{aligned} \mathcal{M}(T) &= \sum_{S \in \text{Conf}_0} \varphi(S) \cdot A(S, T) \\ \varphi(S) &= \sum_{T \in \text{Conf}_0} \mathcal{M}(T) \cdot (-1)^{\#T - \#S} K(T, S). \end{aligned}$$

ii)  $\{\varphi(S)\}_{S \in \text{Conf}_0}$  is a topological basis of  $\mathcal{L}_{\mathbb{A}}$ . That is:

$$\mathcal{L}_{\mathbb{A}} \simeq \varinjlim_{p, q} \left( \prod_{S \in \text{Conf}_0^{pq}} \mathbb{A} \cdot \varphi(S) \right).$$

This means that any  $\mathcal{M} \in \mathcal{L}_{\mathbb{A}}$  is express uniquely as an infinite sum

$$\mathcal{M} = \sum_{S \in \text{Conf}_0^{pq}} \varphi(S) \cdot a_S$$

for some  $p, q \in \mathbb{Z}_{\geq 0}$  and  $a_S \in \mathbb{A}$  for  $S \in \text{Conf}_0^{pq}$ .

(8.2) An explicite formula for  $\varphi(S)$ . For  $S \in \text{Conf}_0$ , let us develop  $\varphi(S) = \prod_{U \in \text{Conf}_0} U \cdot \varphi(U, S)$  for  $\varphi(U, S) \in \mathbb{Q}$ . By a use of (3.4), (5.3), (7.2), one obtain:

$$\varphi(U, S) = \sum_{U=U_1^{k_1} \amalg \cdots \amalg U_m^{k_m}} \frac{(k_1 + \cdots + k_m - 1)! (-1)^{k_1 + \cdots + k_m - 1}}{k_1! \cdots k_m!} \binom{U_1, \dots, U_m}{S}.$$

Here  $(U_1, \dots, U_m)$  means a sequence of  $U_i$  where each  $U_i$  appears with multiplicity  $k_i$ . Let  $\mathcal{M} = \sum_{U \in \text{Conf}_0} U \cdot M_U$  be an element of  $\mathcal{L}_{\mathbb{A}}$ . Then,

$$M_U = \sum_{S \in \text{Conf}_0} \varphi(U, S) \cdot M_S.$$

So combining both formula, one is able to determine all coefficients of  $\mathcal{M}$  from the knowledge of that for  $M_S$  ( $S \in \text{Conf}_0$ ).

*Remark.* In general, an element of  $\mathcal{L}_{\mathbb{A}}$  cannot be expressed by an infinite sum of  $\mathcal{M}(T)$  ( $T \in \text{Conf}_0$ ).

## §9. Grouplike Elements $G_{\mathbb{A}}$

(9.1) Let  $\mathbb{A}$  be an algebra with  $\text{char}=0$  without an idempotent element and  $\mathbb{Q} \subset \mathbb{A}$ . Then one has isomorphisms:  $\exp : \mathcal{L}_{\mathbb{A}} \simeq G_{\mathbb{A}}$  and  $\mathcal{L}_{\mathbb{A}, \text{finite}} \simeq G_{\mathbb{A}, \text{finite}}$ .

(9.2) Generators for  $G_{\mathbb{A}, \text{finite}}$  and  $G_{\mathbb{A}}$ .

**Lemma.** Let  $\mathbb{A}$  be a commutative  $\mathbb{Z}$ -torsionfree algebra with a unit element and without an idempotent element.

i) Any element  $g$  of  $G_{\mathbb{A}, \text{finite}}$  is uniquely expressed as

$$g = \prod_{i \in I} (1 + \mathcal{A}(S_i))^{c_i}$$

for a finite index set  $I$  and  $S_i \in \text{Conf}_0$  and  $c_i \in \mathbb{A}$  for  $i \in I$ . Therefore the correspondence:  $S \mapsto 1 + \mathcal{A}(S)$  induces an isomorphism:

$$\ll \text{Conf} \gg \otimes_{\mathbb{Z}} \mathbb{A} \simeq G_{\mathbb{A}, \text{finite}}$$

where  $\ll \text{Conf} \gg$  is the group generated by the semigroup  $\text{Conf}$ .

ii) The set  $\{\exp(\varphi(S))\}_{S \in \text{Conf}_0} \subset G_{\mathbb{Z}, \text{finite}}$  form a system of free topological generators of  $G_{\mathbb{A}}$ .

## §10. Limit Elements in $\mathcal{L}_{\mathbb{R}}$

We finally introduce the limit element  $\omega_{\Gamma}$  for a finitely generated group  $\Gamma$  with fix generators. A further study of  $\omega_{\Gamma}$  is beyond the scope of the present paper.

(10.1) We equip a classical topology on the  $\mathbb{R}$ -vector space  $\mathcal{L}_{\mathbb{R}}^{pq}$  by:

$$\mathcal{L}_{\mathbb{R}}^{pq} := \varprojlim_n \mathcal{L}_{\mathbb{R}}^{pq} / \overline{\mathcal{J}_n^{pq}} \cap \mathcal{L}_{\mathbb{R}} = \prod_{S \in \text{Conf}_0^{pq}} \mathbb{R} \cdot \varphi(S),$$

the projective limit topology of  $\mathbb{R}$ -vector spaces.  $\mathcal{L}_{\mathbb{R}, \text{finite}}$  is dense in  $\mathcal{L}_{\mathbb{R}}$ . Samely, we put a classical topology on  $\mathbb{R}[[\text{Conf}]]$  by

$$\mathbb{R}[[\text{Conf}^{pq}]] = \varprojlim_n \mathbb{R} \cdot \text{Conf}^{pq} / \mathcal{J}_n^{pq} = \prod_{S \in \text{Conf}^{pq}} \mathbb{R} \cdot S.$$

- i) *The product and coproduct on  $\mathbb{R}[[\text{Conf}^{pq}]]$  are continuous w.r.t. the classical topology.*
- ii) *The classical topology on  $\mathcal{L}_{\mathbb{R}}^{pq}$  is homeomorphic to the topology induced from that on  $\mathbb{R}[[\text{Conf}]]$ .*
- iii) *Let us equip on  $G_{\mathbb{R}}$  a classcial topology induced from  $\mathbb{R}[\text{Conf}]$  as a subspace. Then  $\exp : \mathcal{L}_{\mathbb{R}}^{pq} \rightarrow G_{\mathbb{R}}^{pq}$  is a homeomorphism.*

### (10.2) Limits of equally dividing points.

The correspondence  $1 + \mathcal{A}(T) \in G_{\mathbb{Z}} \mapsto \#T \in \mathbb{Z}$  extends continuously as an additive character from  $G_{\mathbb{R}}$  to  $\mathbb{R}$ , which we denote by  $\mathfrak{X}_{pt}$ . For  $S \in \text{Conf}$ ,  $(1 + \mathcal{A}(S))^{1/\#(S)}$  is called an equally dividing point, where the exponent  $1/\#(S)$  is chosen to get the equality:  $\mathfrak{X}_{pt}((1 + \mathcal{A}(S))^{1/\#(S)}) = 1$ . By a use of logarithm map, we define the set of log-equally dividing points  $\log(\overline{EDP}^{pq}) := \{\mathcal{M}(T)/\#T : T \in \text{Conf}_+^{pq}\}$  in  $\mathcal{L}_{\mathbb{Q}}$  and its closure  $\overline{\log(\overline{EDP}^{pq})}$  w.r.t. the classical topology in  $\mathcal{L}_{\mathbb{R}}$ . The numerical approximation (5.2) will be repeatedly used to show the following Assertion and the next (10.3) Lemma.

*Assertion. 1. The  $\overline{\log(\overline{EDP}^{pq})}$  is a compact convex set.*

2. Let us develop any element  $\omega$  of  $\overline{\log(\overline{EDP}^{pq})}$  as  $\sum_{S \in \text{Conf}_0} \varphi(S) \cdot a_S$ . Then i)  $0 \leq a_S \leq q^{\#S-1}/\#\text{Aut}(S)$  for  $S \in \text{Conf}_0$ , ii) if  $a_S \neq 0$  then  $a_{S'} \neq 0$  for  $S' \leq S$ .

### (10.3) Residual representation of the elements of $\overline{\log(\overline{EDP}^{pq})}$ .

Let  $\omega = \lim_{n \rightarrow \infty} M(T_n)/\#T_n$  be an element of  $\overline{\log(\overline{EDP}^{pq})}$  for a sequence  $\{T_n\}_{n \geq 0}$  of  $\text{Conf}^{pq}$  for some  $p$  and  $q$ . We introduce formal power series

$$\begin{aligned} P(t) &:= \sum_{n=0}^{\infty} \#T_n \cdot t^n \in \mathbb{Z}[[t]], \\ P\mathcal{M}(t) &:= \sum_{n=0}^{\infty} \mathcal{M}(T_n) \cdot t^n \in \mathcal{L}_{\mathbb{Q}}[[t]]. \end{aligned}$$

**Lemma.** Suppose that the radius  $r$  of convergence of  $P(t)$  is positive. Then,

- i)  $P\mathcal{M}(t)$  converges in the same radius  $r$ . That is: for any  $S \in \text{Conf}_0$ ,

$$PM(S, t) := \partial_S P\mathcal{M}(t) = \sum_{n=0}^{\infty} M(S, T_n) \cdot t^n \in \mathbb{Q}[[t]]$$

- converges in the radius  $r$  as for  $P(t)$ .
- ii) The value of the proportion  $PM(t)/P(t)$  converges to  $\omega$  in  $\mathcal{L}_{\mathbb{R}}$  as  $t$  tends to the radius  $r$  along real axis from 0. That is:

$$\omega = \lim_{t \rightarrow r} P\mathcal{M}(t)/P(t) = \sum_{S \in \text{Conf}_0} \varphi(S) \cdot \lim_{t \rightarrow r} PM(S, t)/P(t).$$

iii) If  $P(t)$  and  $PM(S, t)$  extend to meromorphic functions at  $t = r$ , then  $PM(S, t)/P(t)$  is regular at  $t = r$  and

$$\lim_{t \rightarrow \infty} PM(S, t)/P(t) = PM(S, t)/P(t) \Big|_{t=r}.$$

**(10.4) The limit element  $\omega_\Gamma$  for a discrete group  $\Gamma$ .**

Let  $\Gamma$  be a group and let  $G$  be a finite set of generators of  $\Gamma$  such that  $G = G^{-1}$  and  $e \notin G$ . We regard  $\Gamma$  as a Cayley graph. For  $\gamma \in \Gamma$ , put  $\ell(\gamma) := \inf\{n : \gamma = g_1 \cdots g_n \text{ for some } g_i \in G, i = 1, \dots, n\}$ . Define a sequence of increasing graphs  $\Gamma_n := \{\gamma \in \Gamma : \ell(\gamma) \leq n\}$  for  $n \in \mathbb{N}$  and associated generating functions,

$$(1) \quad P_\Gamma(t) := \sum_{n=0}^{\infty} \#\Gamma_n \cdot t^n,$$

$$(2) \quad P_\Gamma \mathcal{M}(t) := \sum_{n=0}^{\infty} \mathcal{M}(\Gamma_n) \cdot t^n.$$

Here  $P_\Gamma(t)$  is well known in literature as growth power series for  $\Gamma$  relative to the generator system  $G$ . We have following facts i) and ii).

- i) For any finitely generated group  $\Gamma$ , the series  $P_\Gamma(t)$  and  $P_\Gamma \mathcal{M}(t)$  converge in a positive radius.
- ii) For a wide class of group  $\Gamma$  with a generator system  $G$ ,  $PM(S, t)$  for  $S \in \text{Conf}_0$  becomes a rational function in  $t$ . This includes examples:
  - a) an automatic group  $\Gamma$  ([2]) such that the regular language  $L(W)$  accepted by the word accepter contains at least one shortest word for every element of the group  $\Gamma$ .
  - b) Hyperbolic or negatively curved groups  $\Gamma$  [5]. Then  $P_\Gamma(t)$  and  $P_\Gamma \mathcal{M}(t)$  are rational functions with a common denominator. (D. Epstein et al [3]).

Owe to i), we introduce the limit element for  $(\Gamma, G)$ :

$$\begin{aligned} \omega_\Gamma &:= \lim_{n \rightarrow \infty} \log(1 + \mathcal{A}(\Gamma_n))^{1/\#\Gamma_n} \\ &= \sum_{S \in \text{Conf}_0} \varphi(S) \cdot \lim_{t \rightarrow r} \frac{P_\Gamma M(S, t)}{P_\Gamma M(pt, t)}. \end{aligned}$$

If  $\Gamma$  belongs to the case ii), owing to (10.3) Lemma iii), the limit process is replaced by a residue calculation of rational functions.

$$\omega_\Gamma = \sum_{S \in \text{Conf}_0} \varphi(S) \cdot \left. \frac{P_\Gamma M(S, t)}{P_\Gamma M(pt, t)} \right|_{t=0} \in \mathcal{L}_{\mathbb{Q}(r)}.$$

Here  $r = 1/\limsup_{n \rightarrow \infty} \sqrt[n]{\#\Gamma_n}$  is the first place of pole of  $P_\Gamma(t)$  on the  $\mathbb{R}_+$ , and is a real algebraic number. We know little about  $\omega_\Gamma$  except that it is defined.

(10.5) *A Kabi-kernel condition on  $\omega_\Gamma$ .* The kabi coefficients induces a map:

$$\sum_{S \in \text{Conf}_0^{pq}} \mathbb{A} \cdot \varphi(S) \simeq \sum_{T \in \text{Conf}_0^{pq}} \mathbb{A} \cdot M(T), \quad \varphi(S) \mapsto \sum_{T \in \text{Conf}_0^{pq}} M(T) \cdot (-1)^{\#T - \#S} K(T, S)$$

for fix  $p$  and  $q \in \mathbb{Z}_{\geq 0}$ . By completing the map, we obtain a map:

$$K : \mathcal{L}_{\mathbb{A}}^{pq} \rightarrow \prod_{T \in \text{Conf}_0^{pq}} \mathbb{A} \cdot M(T).$$

Since  $K$  is bijective on  $\mathcal{L}_{\mathbb{A}, \text{finite}}^{p,q}$ , one gets amazingly  $\ker(K) \cap \log(\overline{EDP^{pq}}) = \phi$ . Let us give a characterization of an element of  $\ker(K) \cap \overline{\log(EDP^{pq})}$ :

*Assertion.* For an element  $\omega = \lim_{n \rightarrow \infty} \mathcal{M}(T_n)/\#T_n$  of  $\overline{\log(EDP^{pq})}$ ,  $\omega$  belongs to  $\ker(K)$ , if and only if  $\lim_{n \rightarrow \infty} \delta(S, T_n)/\#T_n = 0$  for any  $S \in \text{Conf}_0$ , where  $\delta(S, T) := \#\{\text{connected components of } T, \text{ which are isomorphic to } S\}$ .

**Corollary.** The  $\omega_\Gamma$  for an infinite group  $\Gamma$  belongs to  $\ker(K)$ .

## References

1. Baxter, R. J.: Exactly solved models in statistical Mechanics. Academic Press, New York 1982
2. Cannon, J. W., Epstein, D.B.A., Holt, D.F., Paterson, M.S., Thurston, W.P.: Word processing and group theory. Preprint 1988
3. Epstein, D.: A personal communication to the author. March 1990
4. Gromov, M.: Hyperbolic groups, essays in Group theory (edited by S.M. Gersten). MSRI Publications 8. Springer, Berlin Heidelberg New York 1987, pp. 75–263
5. Rota, G.-C.: Coalgebras and Bialgebras in Combinatorics. Lectures at the Umbral Calculus Conference, The University of Oklahoma, May 15–19, 1978. Notes by S.A. Joni
6. Saito, K.: Moduli space for Fuchsian groups. Algebraic Analysis, vol. II (edited by Kashiwara and Kawai). Academic Press, New York 1988, pp. 735–786
7. Saito, K.: The limit element in the configuration algebra for discrete groups. In preparation
8. Sweedler, M.E.: Hopf algebras. Benjamin, 1969

# Some Recent Results on Weakly Pseudoconvex Domains

Nessim Sibony

CNRS URA D 0757, Université de Paris-Sud, Mathématiques, Bâtiment 425  
F-91405 Orsay Cedex, France

We give here a survey of some recent results concerning analysis in weakly pseudoconvex domains in  $\mathbb{C}^n$  with smooth boundary.

Let  $\Omega \Subset \mathbb{C}^n$  be a smoothly bounded domain with defining function  $r$ , more precisely  $r$  is a smooth function defined in a neighborhood  $U$  of  $\bar{\Omega}$  such that  $\Omega = \{z \in U ; r(z) < 0\}$  and  $dr \neq 0$  on  $\partial\Omega$ .

The domain  $\Omega$  is pseudoconvex iff for  $z \in \partial\Omega$  and  $t \in \mathbb{C}^n$

$$\langle Lr(z)t, t \rangle = \sum_{j,k} \frac{\partial^2 r(z)}{\partial z_j \partial \bar{z}_k} t_j \bar{t}_k \geq 0 \quad \text{whenever} \quad \sum_{j=1}^n \frac{\partial r(z)}{\partial z_j} t_j = 0.$$

If the inequality is strict for  $t \neq 0$  the domain is said to be strictly pseudoconvex. Pseudoconvex domains with smooth boundary are just the domains of holomorphy with smooth boundary.

The analysis on strictly pseudoconvex domains received much attention in the late sixties and in the seventies, the explicit construction of kernels in order to solve  $\bar{\partial}$ -equation was one of the main tools to solve function theoretic questions in these domains, see [HeL]. The fact that a strictly pseudoconvex domain is locally biholomorphic to a strictly convex domain is crucial in this approach. Such a simple local model does not exist when the Levi form  $L$  is not positive definite. We give first few examples to show the type of difficulties we have to face.

## 1. Examples

a) The Kohn-Nirenberg example [KN1]. Let  $\Omega = \{(z, w) \in \mathbb{C}^2 ; \operatorname{Re} w + |z|^{2k} + t |z|^2 \operatorname{Re}(z^{2k-2}) < 0\}$ . If  $|t| \leq \frac{k^2}{2k-1}$  then  $\Omega$  is pseudoconvex. If  $|t| > 1$  and  $k \geq 3$  then there is no supporting analytic set to  $\bar{\Omega}$  at the point 0. Hence  $\Omega$  is not biholomorphically equivalent in a neighborhood of 0 to a convex domain.

b) Non embeddability into convex domains. There exists a smooth pseudoconvex domain  $\Omega \Subset \mathbb{C}^3$  and  $p \in \partial\Omega$ , such that for any  $N$  and any convex domain  $U \Subset \mathbb{C}^N$  there is no proper holomorphic map from  $\Omega \cap B(p, r)$ ,  $r > 0$ , into  $U$ , [Si1].

c) The “worm” domain. Diederich and Fornaess [DF1] have exhibited a pseudoconvex domain  $\Omega \Subset \mathbb{C}^2$  with smooth boundary such that  $\bar{\Omega}$  does not have

a Stein neighborhood basis. Moreover given  $\varepsilon > 0$  it is possible to construct a “worm” domain  $\Omega_\varepsilon$  such that there is no plurisubharmonic (p.s.h.) function  $\varrho$  in  $\Omega_\varepsilon$  satisfying  $-Ad^\varepsilon \leq \varrho \leq -Bd^\varepsilon$  with positive constants  $A, B$ . Here  $d$  denote the distance to the boundary of  $\Omega_\varepsilon$ , see [DF1] and [Ki]. This last property implies the non existence of a  $\mathcal{C}^1$  plurisubharmonic defining function for  $\Omega_\varepsilon$ , if  $\varepsilon < 1$ .

Despite the smoothness of  $\partial\Omega$ , the structure of the set  $W(\partial\Omega) = \{z \in \partial\Omega ;$  there exists  $t \in \mathbb{C}^n, t \neq 0, \langle \partial r(z), t \rangle \geq 0, \langle Lr(z)t, t \rangle \geq 0\}$  can be quite wild. Even when  $W(\partial\Omega)$  is reduced to one point, most of the sharp results obtained in the strictly pseudoconvex case either do not generalize or require a more subtle analysis. Depending on the property under consideration one has to introduce specific classes of pseudoconvex domains.

## 2. $B$ -Regular Domains

Let  $U \Subset \mathbb{C}^n$  be a bounded domain. The boundary  $\partial U$  is  $B$ -regular iff for every  $p \in \partial U$  there exists a function  $\psi \in \mathcal{C}(\bar{U})$  p.s.h. in  $U$  with  $\psi(p) = 0$  and  $\psi < 0$  on  $\bar{U} \setminus \{p\}$ . See [Ca1] [Si2] [Si3] for this notion. The following result is proved in [Si2] [Si3].

**Theorem 2.1.** *Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Assume  $\partial\Omega$  is  $B$ -regular. Fix  $0 < \varepsilon < 1$ . There exist two defining functions  $r_1, r_2$  for  $\Omega$  such that  $\varrho_1 = -(-r_1)^{1-\varepsilon}$  is p.s.h. in  $\Omega$  and  $\varrho_2 = (r_2)^{1+\varepsilon}$  is p.s.h. in a neighborhood of  $\Omega$ . Hence  $\bar{\Omega}$  has a Stein neighborhood basis.*

The existence of a Stein neighborhood basis has been proved in many special cases by Diederich-Fornaess [DF2, DF3].

The assumption in the Theorem is the existence of a p.s.h. barrier at every point of the boundary. This is a stronger assumption than the non existence of analytic structure on  $\partial\Omega$ . It can be shown that if such barriers exist except on a “small” set then they exist everywhere [Si2].

When  $W(\partial\Omega)$  is a countable union of closed sets where continuous functions can be approximated by p.s.h. ones, then  $\partial\Omega$  is  $B$ -regular. This is the case when  $\partial\Omega$  is real analytic [DF2] or of finite type [Ca1], in these cases  $W(\partial\Omega)$  is of Hausdorff dimension  $2n - 2$ . However  $B$ -regularity of  $\partial\Omega$  allows the set  $W(\partial\Omega)$  to be of positive Lebesgue measure in  $\partial\Omega$ .

## 3. $\bar{\partial}$ -Problems

Many of the most natural problems that arise for weakly pseudoconvex domains are connected with the question of solving the  $\bar{\partial}$ -equation in the standard spaces  $L^p, H^s, A^s$ .

The most fundamental result in this context is probably the following theorem due to Hörmander [Ho1].

**Theorem 3.1.** *Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain. For every  $g \in L^2_{(0,1)}(\Omega)$  such that  $\bar{\partial}g = 0$  there exists a function  $u$  such that*

- i)  $\bar{\partial}u = g$ ;
- ii)  $\int_{\Omega} |u|^2 \leq C(\text{diam } \Omega) \int_{\Omega} |g|^2$ .

The result holds without any assumption on the smoothness of  $\partial\Omega$ . The question of solving  $\bar{\partial}$ -equation in Sobolev spaces has been settled by Kohn, see [K1].

**Theorem 3.2.** *Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Given  $s_0 > 0$  there exists a constant  $C(s_0, \Omega)$  such that for every closed form  $\alpha \in H_{(0,1)}^s(\Omega)$  with  $s \leq s_0$ , there exists  $u \in H^s(\Omega)$  such that*

$$\bar{\partial}u = \alpha$$

and

$$\|u\|_s \leq C(s_0, \Omega) \|\alpha\|_s.$$

Similar results for the  $\bar{\partial}_b$ -equation have been proved by Kohn [K1], Shaw [S] and Boas-Shaw [BoS].

**Theorem 3.3.** *Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Then for any  $s \geq 0$ ,  $\bar{\partial}_b$  has closed range in  $H_{(0,1)}^s(\partial\Omega)$ .*

Let  $A^s(\Omega)$  denote the space of Hölder continuous functions of order  $s > 0$ , with the usual convention when  $s$  is an integer, for  $s = 0$  we identify  $A^0(\Omega)$  and  $C^0(\bar{\Omega})$ . It is not in general possible to solve the  $\bar{\partial}$ -equation with  $A^s$  estimates as the following result shows, [Si3, Si4].

**Theorem 3.4.** *Let  $m \geq 3$ . There is an  $\Omega \Subset \mathbb{C}^m$ , pseudoconvex with smooth boundary, such that for every  $0 \leq s < \infty$ , there exists a  $\bar{\partial}$ -closed form  $\alpha \in A_{(0,1)}^s(\Omega) \cap H_{(0,1)}^{s+\frac{m-1}{2}}(\Omega)$  such that the equation  $\bar{\partial}u = \alpha$  has no solution in  $A^s(\Omega)$ .*

*Remark.* For  $s$  fixed the domain can be made strictly pseudoconvex except at one point.

Concerning  $L^p$  estimates the problem was studied by Fornaess-Sibony [FSi1]. First of all there is no analogue of Hörmander's result for  $p > 2$ .

**Theorem 3.5.** *There exists a pseudoconvex domain  $U \Subset \mathbb{C}^2$  and a  $\bar{\partial}$ -closed  $(0,1)$  form  $\alpha$ , real analytic in a neighborhood of  $\bar{U}$ , such that the equation  $\bar{\partial}u = \alpha$  has no solution in  $L^p(U)$  for any  $p > 2$ .*

Observe that  $\partial U$  is not smooth. The example is based on the fact that the holomorphic functions in  $L^p(U)$ ,  $p > 2$ , extend holomorphically to a domain  $\tilde{U}$  which is not contained in  $\bar{U}$ . Such phenomenon cannot happen when  $U$  is Runge.

If one tries to solve the  $\bar{\partial}$ -equation in a Hartogs domain  $U = \{(z, w) \in \mathbb{C}^2 ; z \in \Omega \subset \mathbb{C}, |w| < \exp(-\varphi(z))\}$  where  $\varphi$  is subharmonic in an open set  $\Omega \Subset \mathbb{C}$ , one is led to the problem in one complex variable : solve  $\partial u / \partial \bar{z} = f$  in  $L^p(\Omega, \varphi)$  with the estimate

$$\left( \int_{\Omega} |u|^p e^{-\varphi} \right)^{1/p} \leq C \left( \int_{\Omega} |f|^p e^{-\varphi} \right)^{1/p}.$$

When  $1 < p \leq 2$  and  $\Omega \Subset \mathbb{C}^n$  the above estimate holds. The proof is based on the fact that Hörmander's estimates work with the same constant on a domain in  $\mathbb{C}^n$  and on any covering of this domain regardless of the number of sheets.

The situation is different for  $p > 2$ . There exists a subharmonic function  $\varphi$  in the unit disc  $A$  and  $f \in \bigcap_{p>2} L^p(A, \varphi)$  such that for any  $p' > 2$  the equation  $\partial u / \partial \bar{z} = f$  has no solution in  $L^{p'}(A, \varphi)$ . Using these one variable results, we prove the following theorem.

**Theorem 3.6.** *There exists  $\Omega \Subset \mathbb{C}^n$  pseudoconvex with smooth boundary, strictly pseudoconvex except at one point with the following property : for every  $p > 2$  there exists a  $\bar{\partial}$ -closed form  $\beta \in L_{(0,1)}^p(\Omega)$  such that the equation  $\bar{\partial}u = \beta$  has no solution  $u \in L^{p'}(\Omega)$  if  $p' \in ]\sqrt{8+9p}-3, p]$ .*

The result is probably not sharp. There is also a smooth pseudoconvex domain  $\Omega \Subset \mathbb{C}^3$  such that  $\bar{\partial} : L^p(\Omega) \rightarrow L_{(0,1)}^p(\Omega)$  does not have closed range.

The following questions are then quite natural.

i) Suppose  $\Omega \Subset \mathbb{C}^n$  is pseudoconvex with smooth boundary. Let  $s > 0$ , for which  $s' < s$  does the operator  $\bar{\partial} : A^s(\Omega) \rightarrow A_{(0,1)}^s(\Omega)$  has closed range ?

ii) Let  $1 \leq p \leq \infty$ , for which  $p' < p$ , does the operator  $\bar{\partial} : L^{p'}(\Omega) \rightarrow L_{(0,1)}^p(\Omega)$  has closed range ?

The following partial answer due to Bonami-Sibony [BoSi] is a consequence of Kohn's result and of a Sobolev embedding theorem which basically says that if  $\bar{\partial}u \in L^t(\Omega)$  for  $t$  large enough then the dimension  $2n$  in the usual Sobolev embedding theorem may be replaced by  $(n+1)$ .

**Theorem 3.7.** *Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Let  $\alpha \in H_{(0,1)}^s(\Omega)$  a  $\bar{\partial}$ -closed form.*

a) If  $0 < s < \frac{n+1}{2}$  and  $\alpha \in L_{(0,1)}^t(\Omega)$ , with  $\frac{1}{t} = \frac{1}{2} - \frac{s-1/2}{n+1}$ , then the equation  $\bar{\partial}u = \alpha$  has a solution in  $L^r(\Omega)$ , with  $\frac{1}{r} = \frac{1}{2} - \frac{s}{n+1}$ .

b) If  $\frac{n+1}{2} < s \leq \frac{n+2}{2}$  and  $\alpha \in L_{(0,1)}^t(\Omega)$ , with  $\frac{1}{t} = \frac{1}{2} - \frac{s-1/2}{n+1}$ , then  $\bar{\partial}u = \alpha$  has a solution in  $A^\sigma(\Omega)$ , with  $\sigma = s - \frac{n+1}{2}$ .

c) If  $s > \frac{n+2}{2}$  and  $\alpha \in A_{(0,1)}^{\sigma-1/2}(\Omega)$ , then  $\bar{\partial}u = \alpha$  has a solution in  $A^\sigma(\Omega)$ , with  $\sigma = s - \frac{n+1}{2}$ .

*Remarks.* i) In all cases the solution  $u$  belongs to  $H^s(\Omega)$ .

ii) The examples mentioned in Theorem 3.4 show that there is necessarily a loss in the regularity of the solution.

#### 4. $\bar{\partial}$ -Neumann Problem

Let  $\Omega \subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Let  $\alpha$  be a  $\bar{\partial}$ -closed  $(0,1)$  form with  $L^2$  coefficients. An aspect of the  $\bar{\partial}$ -Neumann problem consists in studying the regularity of the solution  $u$  of the equation  $\bar{\partial}u = \alpha$ , where  $u$  is orthogonal to the holomorphic  $L^2$  functions. This solution is called the canonical solution.

Let us say that  $\Omega$  satisfies  $(R_k)$  if the canonical solution  $u$  is in  $H^k(\Omega)$  whenever  $\alpha \in H_{(0,1)}^k(\Omega)$ ;  $\Omega$  satisfies  $(R)$  if the canonical solution  $u$  belongs to  $\mathcal{C}^\infty(\bar{\Omega})$  whenever  $\alpha \in \mathcal{C}^\infty_{(0,1)}(\bar{\Omega})$ .

It is known that property  $(R_k)$  holds for all  $k$ , whenever  $\partial\Omega$  is  $B$ -regular, see Catlin [Ca1]. Very recently, it was proved by Boas-Straube [BoS] that  $(R_k)$  holds for all  $k$  if the domain  $\Omega$  has a p.s.h. defining function.

On the other hand, inspired by ideas of Kiselman [Ki], Barrett has shown in [Ba] that worm domains do not have property  $(R_k)$  for  $k$  large enough, depending on the domain. It seems likely that property  $(R)$  fails in general.

Property  $(R)$  is connected with the boundary behavior of biholomorphic mappings. Indeed Bell-Ligocka have shown, in [BL], that a biholomorphic mapping between two smoothly bounded domains in  $\mathbb{C}^n$ , satisfying property  $(R)$ , extends smoothly up to the boundary.

It is not known whether a biholomorphism between two smoothly bounded domains in  $\mathbb{C}^n$  extends smoothly to the boundary.

#### 5. Domains of Finite Type

A class of domains for which one can hope to generalize the results obtained for the strictly pseudoconvex domains is the class of domains of finite type.

Let  $M$  be a real hypersurface in  $\mathbb{C}^n$  with defining function  $r$ . Let  $p \in M$  and  $\mathcal{V}_p$  the space of germs at  $p$  of irreducible one-dimensional analytic subvarieties. For  $X \in \mathcal{V}_p$  define

$$\tau^*(X) = \sup \left\{ \alpha > 0 : \limsup_{z \in X \setminus \{p\}, z \rightarrow p} \frac{|r(z)|}{\|z - p\|^\alpha} < \infty \right\}.$$

The type of  $M$  at  $p$  is defined as  $\tau(M, p) = \sup\{\tau^*(X), X \in \mathcal{V}_p\}$ .  $M$  is of finite type if  $\sup_{p \in M} \tau(M, p) < \infty$ . This notion has been studied by D'Angelo who showed that  $p \rightarrow \tau(M, p)$  is not uppersemicontinuous but that the set of points of finite type is open, [Dan1], [Dan2]. A point in  $M$  is a point of strict pseudoconvexity iff it is a point of type 2. Bounded pseudoconvex domains with real analytic boundary are of finite type, [DF2].

Let  $U$  be a neighborhood of  $p \in \partial\Omega$ . A subelliptic estimate of order  $\varepsilon$  holds for  $(0,1)$  forms supported in  $U \cap \bar{\Omega}$  iff there exists a constant  $C > 0$  such that

$$\| |u| \|_\varepsilon^2 \leq C(\|\bar{\partial}u\|^2 + \|\bar{\partial}^*u\|^2 + \|u\|^2)$$

for  $u \in \mathcal{D}^{0,1}(U) = \{u ; u \text{ smooth } (0,1)\text{-form supported in } U \cap \bar{\Omega}, \bar{\partial}^*u \in \text{Dom}(\bar{\partial}^*)\}$ , where  $\| |\cdot| \|_\varepsilon$  denotes the tangential Sobolev norm of order  $\varepsilon$ .

Following the work of Kohn [K2] on subelliptic estimates for domains with real analytic boundary, Catlin has proved the following result, [Ca2], [Ca3].

**Theorem 5.1.** Let  $\Omega \Subset \mathbb{C}^n$  be a pseudoconvex domain with smooth boundary. Subelliptic estimate for the  $\bar{\partial}$ -Neumann problem holds for some  $\varepsilon > 0$  in a neighborhood of  $p \in \partial\Omega$  iff  $\partial\Omega$  is of finite type at  $p$ .

The existence of a subelliptic estimate of order  $\varepsilon > 0$  near a point  $p$  implies that the canonical solution to  $\bar{\partial}u = g$  is smooth near  $p$  if  $g$  is smooth in a neighborhood of  $p$ , [KN2].

The best  $\varepsilon$  in the subelliptic estimate is known only for domains of finite type in  $\mathbb{C}^2$ , [FeK, Ca4, FSi2] and for convex domains of finite type in  $\mathbb{C}^n$  [FSi2].

Real progress has been made recently in the understanding of regularity properties of the canonical solution of the  $\bar{\partial}$  and  $\bar{\partial}_b$  equation for domains of finite type in  $\mathbb{C}^2$ , see the survey paper by Christ in these proceedings [Ch].

The understanding of the local geometry of domains of finite type in  $\mathbb{C}^n$  is still incomplete. The following notion of type connected with p.s.h. barriers should be compared to the type  $\tau$ .

For  $p \in \partial\Omega$  let  $\mathcal{P}_p = \{\varphi ; \varphi \text{ p.s.h. continuous in } \bar{\Omega}, \varphi \leq 0, \varphi(p) = 0\}$ . Define for  $\varphi \in \mathcal{P}_p$

$$\tau^+(\varphi) = \limsup_{z \rightarrow p, z \in \Omega} \frac{\log |\varphi(z)|}{\log \|z - \varphi\|}, \quad \tau^-(\varphi) = \liminf_{z \rightarrow p, z \in \Omega} \frac{\log |\varphi(z)|}{\log \|z - p\|}$$

and

$$P(p, \partial\Omega) = \inf_{\varphi \in \mathcal{P}_p} \frac{\tau^+(\varphi)}{\tau^-(\varphi)}.$$

It is shown in [FSi1] that  $\tau(p, \partial\Omega) \leq P(p, \partial\Omega)$ . The converse is proved in [FSi1] when  $\Omega$  is pseudoconvex in  $\mathbb{C}^2$  or convex in  $\mathbb{C}^n$ . The invariant  $P$  is connected with sharp subelliptic estimate.

It would be of interest to compare  $P(p, \partial\Omega)$  and  $\tau(p, \partial\Omega)$  for smooth pseudoconvex domains in  $\mathbb{C}^n$ .

## 6. Spectral Questions and Approximation Problems

For  $0 \leq k \leq \infty$ , let  $A^k(\bar{\Omega})$  denote the algebra of holomorphic functions in  $\mathscr{C}^k(\bar{\Omega})$ . Let  $H^\infty(\Omega)$  denote the algebra of bounded holomorphic functions in  $\Omega$  with uniform norm.

When  $\Omega \Subset \mathbb{C}^n$  is a pseudoconvex domain with smooth boundary, then the Gelfand spectrum of the algebra  $A^k(\bar{\Omega})$ ,  $0 \leq k \leq \infty$ , can be identified with  $\bar{\Omega}$ , [HaSi]. The case of  $H^\infty(\Omega)$  is different. It is not known whether  $\Omega$  is dense in the Gelfand spectrum of  $H^\infty(\Omega)$  even when  $\Omega$  is the unit ball or the unit polydisc in  $\mathbb{C}^n$ ,  $n > 1$ . However there exists a pseudoconvex domain  $\Omega \Subset \mathbb{C}^3$  with smooth boundary, strictly pseudoconvex except at one point of  $\partial\Omega$ , such that  $\Omega$  is not dense in the spectrum of  $H^\infty(\Omega)$ , [Si5].

Very little is known about approximation problems of holomorphic functions defined on a smoothly bounded pseudoconvex domain  $\Omega \Subset \mathbb{C}^n$ . Because of the existence of domains  $\Omega$  such that  $\bar{\Omega}$  does not have a Stein neighborhood basis it is not possible, in general, to approximate uniformly functions in  $A^\infty(\bar{\Omega})$  by holomorphic functions in a neighborhood of  $\bar{\Omega}$ . If we assume that  $\bar{\Omega}$  has a Stein neighborhood basis, it is still not known if functions in  $A^\infty(\bar{\Omega})$  are uniform limits

of functions holomorphic in a neighborhoof of  $\bar{\Omega}$ . This is however true if we assume that  $\partial\Omega$  is  $B$ -regular [Si2].

When  $\Omega \subset \mathbb{C}^n$  is pseudoconvex and smoothly bounded, it is not known if  $A^\infty(\bar{\Omega})$  is dense in  $A^r(\bar{\Omega})$  for the  $C^r$  norm when  $0 \leq r < \infty$ .

For a more detailed discussion of these questions we refer to [Si3].

## References

- [Ba] Barrett, D.: Behavior of the Bergman projection on the Diederich-Fornaess worm. Preprint
- [BL] Bell, S., Ligocka, E.: A simplification and extension of Fefferman's theorem on biholomorphic mappings. *Invent. math.* **57** (1980) 283–289
- [BoS] Boas, H., Shaw, M.C.: Sobolev estimates for the Lewy operator on weakly pseudoconvex boundaries. *Math. Ann.* **274** (1986) 221–231
- [BoSi] Bonami, A., Sibony, N.: Sobolev embedding theorem in  $\mathbb{C}^n$  and the  $\bar{\partial}$ -equation. Preprint
- [BoS] Boas, H., Straube, E.: Sobolev estimates for the  $\bar{\partial}$ -Neumann operator on domains in  $\mathbb{C}^n$  admitting a defining function that is p.s.h. on the boundary. Preprint
- [Ca1] Catlin, D.: Global regularity of  $\bar{\partial}$ -Neumann problem. *Proc. Symp. Pure Math.* **41** (1984) 39–49
- [Ca2] Catlin, D.: Necessary conditions for subellipticity of the  $\bar{\partial}$ -Neumann problem. *Ann. Math.* **117** (1983) 147–171
- [Ca3] Catlin, D.: Subelliptic estimates for the  $\bar{\partial}$ -Neumann problem on pseudoconvex domains. *Ann. Math.* **126** (1987) 131–191
- [Ca4] Catlin, D.: Estimates of invariant metrics on pseudoconvex domains of dimension two. *Math. Z.* **200** (1989) 429–466
- [Ch] Christ, M.: Analysis of  $\bar{\partial}_b$  and  $\bar{\partial}$  on domains of finite type in  $\mathbb{C}^2$ . *These Proceedings*, p. 859
- [Dan1] D'Angelo, J.P.: Subelliptic estimates and failure of semi-continuity of orders of contact. *Duke Math. J.* **47** (1980) 955–957
- [Dan2] D'Angelo, J.P.: Real hypersurfaces, order of contact and applications. *Ann. Math.* **115** (1982) 615–637
- [DF1] Diederich, K., Fornaess, J.E.: An example with non trivial nebenhuelle. *Math. Ann.* **225** (1977) 275–292
- [DF2] Diederich, K., Fornaess, J.E.: Pseudoconvex domains with real analytic boundary. *Ann. Math.* **107** (1978) 371–384
- [DF3] Diederich, K., Fornaess, J.E.: Stein neighborhoods for finite preimage of regular domain. *Manuscripta Math.* **50** (1985) 11–27
- [FeK] Fefferman, C.L., Kohn, J.J.: Hölder estimates on domains of complex dimension two and on three dimensional CR manifolds. *Adv. Math.* **62** (1988) 223–303
- [FSi1] Fornaess, J.E., Sibony, N.: On  $L^p$  estimates for  $\bar{\partial}$ . To appear in *Proc. Symp. Pure Math.*
- [HaSi] Hakim, M., Sibony, N.: Spectre de  $A(\bar{\Omega})$  pour les domaines faiblement pseudoconvexes réguliers. *J. Funct. Anal.* **37** (1980) 127–135
- [HeL] Henkin, G.M., Leiterer, J.: Theory of functions on complex manifolds. Birkhäuser, 1984
- [Ho] Hormander, L.:  $L^2$  estimates and existence theorems for the  $\bar{\partial}$  operator. *Acta Math.* **113** (1965) 89–152
- [Ki] Kiselman, Ch.: A study of the Bergman projection in certain Hartogs domains. Preprint

- [K1] Kohn, J.J.: The range of the tangential Cauchy-Riemann operator. *Duke Math. J.* **53** (1986) 525–545
- [K2] Kohn, J.J.: Subellipticity of the  $\bar{\partial}$ -Neumann problem on pseudoconvex domains: sufficient conditions. *Acta Math.* **142** (1979) 79–122
- [S] Shaw, M.C.:  $L^2$  estimates and existence theorems for the Cauchy Riemann complex. *Invent. math.* **82** (1985) 133–150
- [Si1] Sibony, N.: Sur le plongement des domaines faiblement pseudoconvexes. *Math. Ann.* **273** (1986) 209–214
- [Si2] Sibony, N.: Une classe de domaines pseudoconvexes. *Duke Math. J.* **55** (1987) 299–319
- [Si3] Sibony, N.: Some aspects of weakly pseudoconvex domains. To appear in *Proc. Symp. Pure Math.*
- [Si4] Sibony, N.: On Hölder estimates for  $\bar{\partial}$ . *Ann. Math. Studies*
- [Si5] Sibony, N.: Problème de la couronne pour les domaines pseudoconvexes à bord lisse. *Ann. Math.* **126** (1987) 675–686

# Analysis and Geometry on Groups

Nicholas Th. Varopoulos

Mathématiques, Université de Paris VI, 4 place Jussieu, F-75230 Paris Cedex 05, France

The results that I shall survey here can be seen from several different angles. There is a discrete point of view related to discrete finitely generated groups; there is also a  $C^\infty$  point of view related to connected Lie groups. One can, to a certain extent, unify the above two settings by considering general compactly generated locally compact groups but I shall not do so here. Both in the discrete and in the  $C^\infty$  case we can put forward either the Geometric formulation, such as Sobolev inequalities, or the analytical formulation that examines the behaviour of natural semigroups of operators on  $L^2(G)$ . What makes the theory hold together, in a final analysis, is that equivalence of all these different aspects. To explain how this comes about I have to start with some definitions.

## 1. Distance and Volume Growth

Let  $G$  be a discrete group generated by a finite number of generators  $\gamma_1, \dots, \gamma_k \in G$ . One defines then a distance  $d(\cdot, \cdot)$  on  $G$  by requiring that  $d(gx, gy) = d(x, y)$  ( $x, y, g \in G$ ) and that  $d(e, x)$ , the distance of  $x \in G$  from the neutral point  $e \in G$  is, by definition, the smallest  $n \geq 0$  for which we can write  $x = \gamma_{i_1}^{e_1} \dots \gamma_{i_n}^{e_n}$  ( $i_1, \dots, i_n = 1, \dots, k$ ;  $e_j = 0, \pm 1$ ).

Let  $G$  be a connected Lie group and let  $X_1, \dots, X_k \in \mathcal{L}(G)$  be a finite number of generators of the Lie algebra of  $G$ ; in other words  $X_1, \dots, X_k$  are left invariant  $C^\infty$  fields on  $G$  that together with their successive brackets  $[X_{i_1}[X_{i_2}, \dots, X_{i_s}], \dots]$  generate the tangent space. We say that an absolutely continuous path  $\dot{l}(t) \in G$  ( $0 \leq t \leq T$ ) is of length less or equal to  $T$  if its speed vector  $\dot{l}(t) = dl\left(\frac{\partial}{\partial t}\right)$  (with respect to  $X_1, \dots, X_k$ ) is almost everywhere of length  $\leq 1$ : This means that  $\dot{l}(t) = \sum_{j=1}^k a_j X_j$  (p.p.  $t \sum a_j^2 \leq 1$ ). We then say that  $d(x, y) \leq T$  ( $x, y \in G$ ) if we can join  $x$  to  $y$  with a path of length  $\leq T$ .

The growth function  $\gamma(t)$  ( $t > 0$ ) of  $G$  is in either of the above two cases defined to be  $\gamma(t) =$  The Haar measure of a ball of radius  $t$ . For large  $t$  ( $t \geq 1$ ) the above function  $\gamma(t)$  is essentially independent of the particular choice of the generators used:  $\gamma(t)$  ( $t \geq 1$ ) is thus a group invariant. For Lie groups and  $0 < t < 1$  the behaviour of  $\gamma(t)$  does depend on the choice of  $X_1, \dots, X_k$  but we always have

$\gamma(t) \approx t^\delta$  with  $\delta = \delta(G, X_1, X_2, \dots, X_k) = 1, 2, \dots$ . (This is a theorem of Nagel-Stein-Wainger). For  $t \geq 1$  and a Lie group we have either  $\gamma(t) \approx t^D$  with  $D = D(G) = 0, 1, \dots$  or  $\gamma(t) \geq Ce^{ct}$  (This is a theorem of Guivarc'h). In the discrete case we have either  $\gamma(t) \approx t^D$  ( $D(G) = 0, 1, 2, \dots$ ) if  $G$  is a finite extension of a nilpotent group or  $\gamma(t)t^{-A} \rightarrow \infty$  for all  $A \geq 1$  in all other cases (this is a theorem of Gromov).

## 2. The Diffusion and the Random Walks

Let  $G$  be a unimodular Lie group with  $X_1, \dots, X_k \in \mathcal{L}(G)$  as above, we can then consider  $\Delta = -\sum X_j^2$  which can be identified with a self adjoint (positive) operator on  $L^2(G)$  and we can also consider  $T_t = \exp(-t\Delta)$  the corresponding submarkovian semigroup. The kernel of that semigroup will be denoted by  $p_t(x, y) = p_t(x^{-1}y)$  ( $t > 0; x, y \in G$ ). The discrete analogue of the above diffusion is of course the random walk defined on a discrete group by the transition matrix  $M(x, y) = \mu(x^{-1}y)$  ( $x, y \in G$ ) where  $\mu \in \mathbb{P}(G)$  is a symmetric probability measure on  $G$ . We shall consider in what follows, essentially, only random walks that are defined by symmetric measures that have generating supports ( $Gp(\text{supp } \mu) = G$ ). What we shall examine then is the convolution powers  $\mu^n$  of that measure or equivalently  $T_t = \exp(-t(\delta - \mu))$  the continuous time Markov semigroup that it generates.

## 3. Analytic and Probabilistic Formulation

One of the main accomplishments of the present methods is that it allows us to study the convolution powers of a finitely supported symmetric measure as considered in the previous section.

**Theorem 1.** *Let  $G$  be a discrete finitely generated group and let  $\mu$  be a measure as above. Let us also assume that  $\gamma(t) \geq ct^D$  for some  $c, D > 0$ . We then have  $\mu^n(\{e\}) = O(n^{-D/2})$ .*

The above theorem allows us in particular to classify the discrete groups for which the series  $\sum \mu^n(\{e\}) = +\infty$ . Such groups are called recurrent groups, the reason being that the random walk with transition matrix  $M(x, y) = \mu(x^{-1}y)$  is a recurrent random walk (and this fact is independent of the particular choice of  $\mu$ ):

**Corollary.** *The only recurrent groups are the finite extensions of the following three groups:  $\{0\}, \mathbb{Z}, \mathbb{Z}^2$ .*

Theorem 1 easily generalizes to convolution products  $\mu_1 * \dots * \mu_n$  provided of course that the measures  $\mu_j$  satisfy the appropriate conditions uniformly in  $j$ . Theorem 1 is a typical result of the discrete version of our theory. The continuous variant of the same result is the following.

**Theorem 2.** Let  $G$  be a unimodular Lie group and let  $X_1, \dots, X_k \in \mathcal{L}(G)$  be as before. Let us assume that the induced growth function satisfies  $\gamma(t) \approx t^\delta (t \rightarrow 0)$  and  $\gamma(t) \geq ct^D$  ( $t \geq 1$ ) for some  $\delta, D = 0, 1, \dots$ .

We then have  $\|p_t\|_\infty = O(t^{-\delta/2})$  ( $t \rightarrow 0$ ) and  $\|p_t\|_\infty = O(t^{-D/2})$  ( $t \rightarrow \infty$ ).

The small time behaviour of  $\|p_t\|_\infty$  described in the Theorem is contained in a previous more general result of A. Sanchez-Calle. The group structure is not essential for this small time behaviour of  $p_t$ . The above two theorems can of course be unified to a single result on locally compact groups and the methods of the proofs, as we shall see, have very little to do with “real analysis”.

The metric  $ds^2 = \varphi(y)(dx^2 + dy^2)$  on  $\mathbb{R}^2$  where  $\varphi(y) = y^{-2}$  for  $|y| \geq 1$  gives an example of a Riemannian manifold that has exponential volume growth (since for  $|y| \geq 1$  it is just the hyperbolic plane) but has “slow” decay for its canonical  $p_t$  as  $t \rightarrow \infty$ . Indeed the above metric is conformal with the Euclidean metric and therefore has no Green’s function i.e.  $\int_1^\infty p_t = +\infty$ . This shows that the group structure in Theorem 2 is essential for the behaviour of  $p_t$  as  $t \rightarrow \infty$ .

## 4. Geometric Formulation

Let  $G$  be a unimodular Lie group and let  $X_1, \dots, X_k \in \mathcal{L}(G)$  be as before. We shall denote the corresponding gradient by:  $\nabla f = (X_1 f, \dots, X_k f) \in \mathbb{R}^k$  ( $f \in C_0^\infty(G)$ ). The main Geometric Theorem is

**Theorem 3.** Let  $G$  and  $X_1, \dots, X_k$  be as before and let  $\delta, D \geq 0$  be as in Theorem 2. Let also  $n \geq 1, \delta \leq n \leq D$  we then have

$$\|f\|_{n/(n-1)} \leq C \|\nabla f\|_1; \quad f \in C_0^\infty.$$

Conversely if the above Sobolev inequality holds for some  $n$  then  $n \geq \delta$  and  $\gamma(t) \geq ct^n$  ( $t \geq 1$ ).

(All the  $\| \cdot \|_p$  norms in what follows are taken in  $L^p(G)$  for the Haar measure).

The above Sobolev type estimates are usually reformulated by the Geometers in terms of isoperimetric inequalities of the type  $|A|_r^{(n-1)/n} \leq C|\partial A|_{r-1}$  ( $A \subset G$ ) where  $|\cdot|_s$  refers to the appropriate  $s$ -dimensional Hausdorff measure and  $r$  is the topological dimension of  $G$ . The discrete analogue of the above theorem states:

**Theorem 4.** Let  $G$  be a discrete finitely generated group, then the Sobolev inequality  $\|f\|_{n/(n-1)} \leq C \|\nabla f\|_1$  ( $f \in C_0(G)$ ) holds for some  $1 \leq n \in \mathbb{R}$  if and only if  $\gamma(t) \geq ct^n$  ( $t \geq 1$ ).

In the above theorem the  $L_1$ -norm of the gradient is of course  $\|\nabla f\|_1 = \sum_{d(x,y)=1} |f(x) - f(y)|$ . Once more the above result can be stated in terms of discrete isoperimetric inequalities.

## 5. The Connection Between Analysis and Geometry

What unifies the Geometric and the Analytic point of view and what is, in a final analysis, the pivot of the proofs is the following general result of Functional Analysis:

Let  $A$  be the generator of an appropriate semigroup  $T_t = e^{-tA}$  of operators on  $L^p(X; dx)$  (the spaces of  $p$ -integrable functions on an abstract measure space  $(X; dx)$ ). Let  $n > 2$ , the following two conditions are then equivalent:

- (i)  $\|f\|_{2n/(n-2)} \leq C(Af, f)^{1/2}; f \in \text{Dom}(A).$
- (ii)  $\|T_tf\|_\infty \leq ct^{-n/2} \|f\|_1; t > 0, f \in L^1.$

For the semigroups associated to our random walks on discrete groups the generator is:  $A = \delta - \mu$  and the Dirichlet form satisfies  $D_\mu(f) = (Af, f) \approx D_0(f)$  where we denote by  $D_0(f) = \sum_{d(x,y)=1} |f(x) - f(y)|^2$  the “standard” Dirichlet form on  $G$ . This equivalence  $D_\mu \approx D_0$  is trivial to see if  $\mu$  has *finite* support but what is important is that it remains true for a more general class of measures; namely for all symmetric Probability measures on  $G$  with generating support and whose “variance” is finite:

$$E(\mu) = \sum_{x \in G} d^2(e, x)\mu(\{x\}) < +\infty.$$

This observation although not very difficult to prove is absolutely crucial for us.

In the case of a Lie group the Dirichlet form of our semigroup  $T_t = e^{-tA}$  is of course the familiar expression

$$(Af, f) = \|Vf\|_2^2 = \int_G \sum |X_j f|^2.$$

Observe finally that the  $L^1 \rightarrow L^\infty$  operator norm  $\|e^{-tA}\|_{1,\infty}$  on a Lie group is  $p_t(e)$  and similarly  $\|e^{-n(\delta-\mu)}\|_{1,\infty} \sim \mu^n(e)$  for a discrete group (This last  $\sim$  has to be interpreted correctly but it certainly implies  $\mu^n(e) = O(n^{-\alpha}) \Leftrightarrow \|e^{-t(\delta-\mu)}\|_{1,\infty} = O(t^{-\alpha})$ ).

With the above facts in mind the connection between the Geometric and the Analytic theory becomes obvious. Another thing that becomes apparent (and this is the single most important feature of all the proofs) is that *changing the measure*  $\mu$ , say in Theorem 1, *makes no difference* as long as we restrict ourselves to measures of finite variance. Indeed such changes leave invariant (up to equivalence) the Dirichlet form  $D_\mu(f)$ . What remains to be done to complete the proof of, say Theorem 1, is to produce *one* symmetric probability measure with finite variance and with convolution powers that decay optimally:  $\mu^n(\{e\}) = O(n^{-D/2})$ .

This last step is done “by hand”. We simply try out a measure of the form:  $\mu = \sum \lambda_j \chi_j$  where  $\lambda_j \geq 0, \sum \lambda_j = 1$  and where  $\chi_j$  denotes the normalized characteristic function of the  $j$ -ball in  $G$ . The condition  $E(\mu) < +\infty$  is easy to express in terms of the  $\lambda$ 's and the convolution powers  $\mu^n$  can be estimated by an elementary argument. The above construction does not seem to work if we restrict ourselves to measures of finite support and this is something that to this day I cannot really explain to myself in a satisfactory manner.

## 6. Further Development and Open Problems

In the interaction between the Geometric and Analytical results there is one point that remains obscure. Indeed what is natural to consider from the semigroup point of view is the norm  $\|\mathcal{A}^{1/2}f\|_p$ , ( $f \in C_0^\infty$ ) ( $\mathcal{A}$  is a self adjoint positive operator) and what occurs naturally in the Geometric formulation is the norm  $\|\nabla f\|_p$ . It is only for  $p = 2$  that the two norms are obviously equivalent and it is an open problem whether we have in general  $\|\nabla f\|_p \approx \|\mathcal{A}^{1/2}f\|_p$ . That this is the case in the real variable situation  $G = \mathbb{R}^n$  is the content of the classical M. Riesz theorem (for  $p \neq 1, \infty$ ). This equivalence holds when  $G$  is a group of polynomial growth (This is a recent theorem of G. Alexopoulos). It also holds when  $G$  is non amenable e.g. a classical non compact semi-simple group (this is a result of N. Lohoué). The problem for a general unimodular group remains open and seems difficult. The above problem has an obvious discrete formulation that contains, no doubt, the essence of the difficulty.

When the group  $G$  is not unimodular then, as we already pointed out, the geometric aspect of our theory goes through in a very satisfactory fashion. What remains very much open is the analytical theory. Indeed the long time behaviour of the appropriate heat kernel remains untractable by the above methods. The problem is very much connected with the analysis of the canonical heat kernel on symmetric spaces. Indeed any symmetric space of non-compact type can be realized, by the Iwasawa decomposition KAN, as the *non-unimodular* group AN.

The last problem that I shall consider consists in obtaining a finer analysis of the behaviour of  $p_t$ , as  $t \rightarrow \infty$ , for Lie group, or  $\mu^n$  for a discrete group. Assume that  $G$  is a unimodular Lie group. If  $G$  is not amenable, and only then, we have  $p_t(e) = O(e^{-\lambda t})$  where  $\lambda > 0$  is the spectral gap of  $\mathcal{A}$  and depends on the particular choice of the fields  $X_1, \dots, X_k$ . There are good reasons to suspect that in fact  $p_t(e) \sim t^{a/2}e^{-\lambda t}$  where  $a$  is some integer or possibly " $+\infty$ " that only depends on the group and not on the choice of the fields (just as for amenable groups where we have  $\lambda = 0$ ). The analogous conjecture for discrete groups is false (the counter example is due to D. Cartwright). If  $G$  is semi-simple this is, once more, related to the heat kernel on symmetric spaces (Ph. Bougerol has examined this case).

For a Lie group of polynomial growth G. Alexopoulos has proved a "local Central Limit" theorem:  $p_t(e)t^{D/2} \xrightarrow[t \rightarrow \infty]{} \alpha_0 > 0$ . The following asymptotic development  $p_t(e) \sim t^{-D/2}[\alpha_0 + \alpha_1 t^{-1/2} + \dots]$  should hold, but this is an open problem. Similarly for semi-simple groups and symmetric spaces C. Herz conjectures that  $p_t(e) \sim e^{-\lambda t}t^{a/2}[\alpha_0 + \alpha_1 t^{-1/2} + \dots]$  (as  $t \rightarrow \infty$ ). Some logarithms could possibly appear in these asymptotic developments.

Let  $G$  be a discrete group and let  $\mu, \nu \in \mathbb{P}(G)$  be two symmetric probability measure of finite variance. Let us also assume that  $\mu^n(e) = O[\exp(-\alpha(n))]$  where  $\alpha(t) \geq 0$  is an increasing positive function of ( $t \geq 0$ ). By a slight variance of the previous methods (here we make essential use of E.B. Davies work in the subject) we can then show that:  $\nu^n(e) = 0 [\exp(-c\tilde{\alpha}(cn))]$  for some  $0 < c$ , where we denote by

$$\tilde{\alpha}(t) = \frac{1}{t} \int_0^t \alpha(s) ds.$$

The analogous result for unimodular Lie groups also holds. What makes this fact interesting is that for many natural functions e.g.  $\alpha(t) = t^\alpha$ ,  $t^\alpha \log(1+t)$  e.c.t. we have  $\alpha \approx \tilde{\alpha}$ . This fact is used to analyse the groups that have *superpolynomial* growth:

Assume that  $G$  (discrete or Lie, amenable or not, but unimodular) satisfies the growth condition  $\gamma(t) \geq \exp(ct^\alpha)$ ,  $t \geq 1$ , for some  $0 < \alpha \leq 1$  [cf. R. Grigorchuk's paper in these proceedings]. Using our methods then we can easily establish that, say for a discrete group, we have  $\gamma(n+1) - \gamma(n) \geq \exp(cn^\alpha)$  with possibly a different  $c > 0$  but the same  $0 < \alpha \leq 1$ . Using this fact and refining our methods further (we use in particular here an idea of L. Saloff-Coste) we can then prove that (again for a discrete group) we have:

$$\mu^n(\{e\}) = O[\exp(-cn^{\alpha/(\alpha+2)})]$$

The analogous result when  $G$  is a Lie group and  $\alpha = 1$  also holds. The above estimate is optimal. Indeed for any non virtually Nilpotent polycyclic group and every finitely supported symmetric  $\mu \in \mathbb{P}(G)$  we have  $\mu^{2n}(e) \geq C \exp[-cn^{1/3}]$  (this was shown by G. Alexopoulos) and for all these groups  $\alpha = 1$ . The details of the above result will appear elsewhere.

A decay of the type  $\exp(-cn^\beta)$  for  $p_t(e)$  gives rise of course to Orlicz type Sobolev inequalities of the form  $\|f\|_{L \log^\gamma L} \leq C \|\nabla f\|_1$  where  $\gamma = \gamma(\beta)$ . In terms of isoperimetric inequalities for discrete groups for instance, we can say that if  $\mu^n = O[\exp(-cn^{-\beta})]$  ( $0 < \beta \leq 1$ ) then we have:

$$|\partial\Omega| \geq C|\Omega|(\log |\Omega|)^{\frac{\beta-1}{\beta}}$$

for all finite  $\Omega \subset G$  with  $|\Omega| \geq 2$  where  $|\cdot|$  denotes the cardinality of a finite set.

For exponential groups this gives:

$$|\partial\Omega| \geq C|\Omega|(\log |\Omega|)^{-2}.$$

A final result that I shall mention concerns  $p_t(x, y)$  the canonical heat kernel on a Riemannian manifold that covers normally some compact manifold with deck transformation group  $G$ . With the present methods we can show that the behaviour of  $\|p_t\|_\infty$  (as  $t \rightarrow \infty$ ) is "identical" with the behaviour of  $\mu^n(e)$  for  $\mu \in \mathbb{P}(G)$  (as in Section 3). The term "identical" means for instance that  $\mu^n(e) = O(n^{-\alpha}) \Leftrightarrow \|p_t\|_\infty = O(t^{-\alpha})$  or more generally that we have the:

$$O[\exp(-\alpha(\cdot))] \Leftrightarrow O[\exp(-c\tilde{\alpha}(c))]$$

correspondence that we considered above. This is one of the very first results that I obtained in the subject and it is this that convinced me of the fundamental connection that existed between the discrete and the continuous theory.

In this survey I have said nothing about the Gaussian estimates of the heat kernels. It would take a different paper to do that. The interested reader could consult the literature below.

## Literature

The theory that we surveyed in this paper is the subject matter of a forthcoming book [1]. A preliminary version of this book exists in the form of mimeographed notes: University Paris VI.

Most of the results that I presented were developed by the author in a series of papers the most significant being [2].

For the Functional Analytic tools of Section 5, and the work of E.B. Davies cf. [3, 4, 5, 6, 12].

For further developments in locally compact groups, cf. [7, 8, 9, 13]. Most of the work of G. Alexopoulos has not yet appeared in print, cf. [10, 14]. For the Symmetric space point of view, cf. [11, 15].

1. N. Th. Varopoulos, L. Saloff-Coste, Th. Coulhon: Analysis and geometry on groups. Cambridge University Press (to appear)
2. N. Th. Varopoulos: Analysis on Lie groups. *J. Funct. Anal.* **76**, no. 2 (1988) 346–410
3. N. Th. Varopoulos: Hardy-Littlewood theory for semigroups. *J. Funct. Anal.* **63**, no. 2 (1985) 240–260
4. P. Bénilan: Operateurs accrétiifs et semi-groupes dans les espaces  $L^p$  ( $1 \leq p \leq +\infty$ ). In: Functional analysis and numerical analysis. Japan-France Seminar
5. A. Yoshikawa: Fractional powers of operators, interpolation theory and imbedding theorems. *J. Fac. Sci. Univ. Tokyo, I.A.* **18** (1971)
6. Th. Coulhon: Semigroup theory and evolution equations, Clement, Mitidien, De Pagter, ed. Marcel Dekker 1991
7. L. Saloff-Coste: Inégalité de Sobolev produite sur les groupes de Lie nilpotents. *J. Funct. Anal.* **79** (1) (1988) 44–56
8. L. Saloff-Coste: Sur la décroissance des puissances de convolution sur les groupes. *Bull. Sci. Math., 2<sup>e</sup> série* **113** (1989) 3–21
9. L. Saloff-Coste: Analyse sur les groupes de Lie à croissance polynomiale. *Arkiv for Matematik* **28** (2) 1990 315–331
10. G. Alexopoulos: *C.R. Acad. Sci. Paris* **309** (I) (1989) 661–662 and **305** (I) (1987) 777–779
11. N. Lohoué: Estimées de type Hardy pour l'opérateur  $A + \lambda$  d'un espace symétrie de type non compact. *C.R. Acad. Sci. Paris* **308** (I) (1989) 11–14
12. E.B. Davies: Heat kernels and spectral theory. Cambridge University Press, 1990
13. N.Th. Varopoulos: ICM-90 Satellite Conference Proceedings Harmonic Analysis (Sendai 1990)
14. G. Alexopoulos: ICM-90 Satellite Conference Proceedings Harmonic Analysis (Sendai 1990) and *Canadian J. Math.* (to appear)
15. P. Bougerol: *Ann. Sci. Ec. Norm. Sup. 4<sup>e</sup> série* **14** (1981) 403–432



# Asymptotically Holomorphic Functions and Certain of Their Applications

Alexander L. Volberg

University of Kentucky, Lexington, KY 40506-0027, USA

## 1. Introduction

In one-dimensional complex analysis the essential part is played by approximation results and by dual uniqueness theorems. One of the general problems of rational approximation may be stated as follows. Let  $X$  be a linear topological space of functions,  $\Lambda$  be a subset of  $\mathbb{C}$  such that fractions  $z \rightarrow (z - \lambda)^{-1}$ ,  $\lambda \in \Lambda$ , lies in  $X$ . How one can find out what subsets  $\Lambda'$ ,  $\Lambda' \subset \Lambda$ , will bring us a complete in  $X$  family of rational fractions  $\{(z - \lambda)^{-1}\}_{\lambda \in \Lambda'}$ ? Quite often it turns out that functions  $f_{x^*}(\lambda) \stackrel{\text{def}}{=} \langle x^*, (z - \lambda)^{-1} \rangle$ ,  $\lambda \in \Lambda$ ,  $x^* \in X^*$ , are “holomorphic” in  $\Lambda$  in a certain (sometimes pretty weak) sense. Thus we are dealing with a description of zero-sets of the class  $\{f_{x^*}\}_{x^* \in X^*}$  of “holomorphic” functions. According to this it seems important to understand the nature of this “holomorphicity.” The appropriate notion may be grasped by the conception of asymptotically holomorphic (AH) functions. These are (very roughly speaking) the functions defined in subdomains of  $\mathbb{C}$  with vanishing  $\bar{\partial}$ -derivative on the boundary. More precisely (and we will see it below) only functions with some critical rate of decreasing of  $\bar{\partial}$  deserved to be called AH.

To begin with we state here two typical results of AH function theory. Then we are going to show how to apply these results to

1. Problems in harmonic analysis (ideal description in the spirit of Domar);
2. A problem in dynamical systems (finiteness of number of limit cycles for QA vector fields);
3. Problems in approximation theory.

Perhaps, it is worthwhile to mention that 3 is the place, where the AH technique appeared. In what follows  $\mathbb{D}$  is the unit disc,  $\mathbb{C}^+$  is the right half-plane.

**Theorem 1.** *Let  $w$  be a positive increasing function on  $(0, 1)$  such that*

$$x \log \frac{1}{w(x)} \uparrow \infty, \quad x \downarrow 0; \quad (\text{R1})$$

$$\int_0^1 \log \log \frac{1}{w(x)} dx = \infty. \quad (1)$$

And let  $f \in C^1(\mathbb{D}) \cap L^\infty(\mathbb{D})$  and

$$|\bar{\partial}f(z)| \leq w(1 - |z|). \quad (2)$$

Then either there exist constants  $c_1, c_2$  such that

$$|f(x)| \leq c_1 w(c_2(1 - |z|)), \quad z \in \mathbb{D} \quad (\text{S1})$$

or

$$\int_{\partial\mathbb{D}} \log |f| dm > -\infty. \quad (\text{B1})$$

**Theorem 1.2.** Let  $f \in C^1(\mathbb{C}^+) \cap L^\infty(\mathbb{C}^+)$  be such that

$$|\bar{\partial}f(z)| \leq w\left(\frac{1}{\operatorname{Re} z}\right), \quad (3)$$

where  $w$  is a positive increasing function on  $(0, 1)$  which satisfies

$$x^2 \frac{\log \frac{1}{w(x)}}{\log \frac{1}{w(2x)}} \uparrow \infty, \quad x \downarrow 0. \quad (\text{R2})$$

Then either there exist constants  $c_1, c_2$  such that

$$|f(x)| \leq c_1 w\left(\frac{c_2}{\sqrt{x}}\right) \quad (\text{S2})$$

or there exists  $n$  such that

$$\lim_{x \rightarrow +\infty} |f(x)| e^{nx} > 0. \quad (\text{B2})$$

**Theorem 1.3.** Let  $w$  be a positive increasing function on  $(0, 1)$  such that for a positive  $\varepsilon$

$$x^\varepsilon \log \frac{1}{w(x)} \uparrow \infty, \quad x \downarrow 0. \quad (\text{R3})$$

And let  $f \in C^1(\mathbb{D} \setminus \{0\})$  and have estimates

$$|f(z)| \leq \left(\frac{1}{w(|z|)}\right)^{\frac{\varepsilon}{\log 1/\varepsilon}} \quad (4)$$

and

$$|\bar{\partial}f(z)| \leq w(|z|). \quad (5)$$

Then either there exist constants  $c_1$  and  $c_2$  such that

$$|f(z)| \leq c_1 w(c_2|z|) \quad (\text{S3})$$

or there exists a function  $f^*$  holomorphic in a punctured neighbourhood of the origin such that

$$\forall c_2 \exists c_1 \quad |f(z)| > |f^*(z)| - c_1 w(c_2|z|)).$$

Now let us discuss these results. These are theorems of “forced behaviour” of AH functions. They serve as a kind of the substitute for the uniqueness theorems for holomorphic function. In fact, a typical uniqueness theorem tells us that a holomorphic function either vanishes identically or is quite big. For AH functions we have a similar alternative: an AH function is either as small as its  $\bar{\partial}$ -derivative ((S1), (S2), (S3)) or is substantially big ((B1), (B2), (B3)). To obtain such a dichotomy one obviously has to have  $w$  in  $|\partial f| \leq w$  fast decreasing, which is indicated in conditions (1), (R2), (R3). The rate of this decreasing is sharp only in Theorem 1.1 and here we must impose an auxiliary conditions (R1) which is a regularity condition. (By the way, (R1) is indispensable in a sense.) (R2), (R3) are certainly too strong. But these actual versions of Theorems 1.2 and 1.3 have the advantage that regularity conditions and “rate-of-decay” conditions on  $w$  are united.

Assumption (4) deserves a separate comment. It singles out Theorem 1.3 from the row. Actually this is not the case. We could state Theorems 1.1, 1.2 in more general form, namely with  $f$  not bounded but with a mild growth (toward  $\partial\Omega$  in Theorem 1.1 and towards  $\infty$  in Theorem 1.2). The reason why we chose these wordings is that we will use results only in this form. Also it is worthwhile to mention that Theorem 1.3 is absolutely trivial in the case of bounded  $f$ . As for Theorems 1.1, 1.2 all technical difficulties are here even for the case of bounded  $f$ .

## 2. Harmonic Analysis Applications

Our main object here is a weighted  $l^2$  space, namely

$$l^2(p) \stackrel{\text{def}}{=} \{(a_n)_{n \in \mathbb{Z}} : \sum_{n \in \mathbb{Z}} |a_n|^2 e^{-p_n} < \infty\},$$

and the right-shift operator  $\tau$ ,

$$\tau : \{a_n\} \rightarrow \{a_{n-1}\}.$$

We will consider only weights with some regularity, and, for instance,

$$\lim_{n \rightarrow \infty} \frac{|p_n|}{\log |n|} = \infty. \quad (6)$$

Moreover we will say that  $p$  is of concave type if  $n \rightarrow p_n$ ,  $n > 0$ , is concave and nonnegative,  $n \rightarrow p_n$ ,  $n < 0$ , is convex negative. We will say that  $p$  is of convex type if  $n \rightarrow p_n$ ,  $n > 0$ , is convex and nonnegative,  $n \rightarrow p_n$ ,  $n < 0$ , is concave and negative. Our weights will be only either of concave or convex type. The description of  $\tau$ -invariant subspaces in  $l^2(p)$  is one of the model problems in harmonic analysis. Substantial contribution is made by Y. Domar [1], [2], [3] and A. Borichev [4], [5]. The idea of using AH function to these problems is due to A. Borichev. We follow his ideas here. This problem is still not fully understood

in spite of considerable efforts undertaken since the classical Beurling description in the case

$$p_n = \begin{cases} 0, & n > 0 \\ -\infty, & n < 0. \end{cases} \quad (7)$$

It is natural to expect that if  $\{p_n\}$  is pretty close to (7) then the structure of invariant subspace lattice is also classical. To state one positive result in this direction let us consider the sequence  $p = (p_n)$  of concave type and such that

$$\left. \begin{array}{l} p_n = 0, \quad n > 0; \\ \frac{|p_{-n}|}{n^{1/2}} \uparrow \infty, \quad n \rightarrow +\infty; \\ \sum_{n \geq 1} \frac{|p_{-n}|}{n^2} = \infty. \end{array} \right\} \quad (8)$$

Let  $c(p) \stackrel{\text{def}}{=} \bigcup_{c>0} l^2(cp)$  be the projective limit of spaces  $l^2(cp)$ ,  $c > 0$ .

**Theorem 2.1.** *Under the assumptions (8) all  $c(p)$ -closed invariant subspaces  $E$  of  $\tau$  such that  $\tau E \neq E$  are of the type  $IH^2$ , where  $I$  is a unimodular function,  $I \in c(p)$ .*

The Hilbert space setting is supposed to be much more difficult. For  $l^2(p)$  with  $p_n = 0$ ,  $n \geq 0$ ;  $p_n = e^{-cn}$ ,  $c > 0$ ,  $n < 0$ , the investigation of  $\tau$ -invariant subspaces was started in [6] and finished in [7].

Remind that for a given space of sequences  $S$  is a subspace ( $k \in \mathbb{Z}$ )

$$S_k = \{\{s_m\} \in S : s_m = 0, m < k\}$$

is called standard. Any standard subspace is  $\tau$ -invariant. When is the converse true? This seems to be interesting as operators with small amount of invariant subspaces naturally attract attention.

We need the result that follows for illustration. This result is weaker than similar results in [1], [2].

Let  $p = (p_n)$ ,  $p_{-n} = -p_n$ , be of convex type and

$$\left. \begin{array}{l} \lim_{n \rightarrow \infty} \frac{p_n}{n \log n} = \infty \\ \lim_{n \rightarrow \infty} (u_{\pm n} - u_{\pm 2^n}) > 0, \quad \text{where } u_n = p_n - p_{n+1} \end{array} \right\} \quad (9)$$

and  $l(p) = \{(a_n) : \forall c_1 > 0 \sum_{n<0} |a_n|^2 e^{-c_1 p_n} < \infty, \exists c_2 > 0 \sum_{n>0} |a_n|^2 e^{-c_2 p_n} < \infty\}$ .

**Theorem 2.2.** *Under the assumption (9) all  $l(p)$ -closed  $\tau$ -invariant subspaces  $E$  of  $l(p)$  such that  $\tau E \neq E$  are standard.*

The proofs are based on results of the first section. We outline briefly these proofs. In both theorems we need to solve convolution equation

$$(a * b)_n = 0, \quad n \leq 0, \quad (10)$$

$a \in c(p)$  (or  $l(p)$ ),  $b \in c(p)^*$  (or  $l(p)^* = l(p)$ )

The usual Fourier transform  $(a) \rightarrow \sum_{n \in \mathbb{Z}} a_n z^n$  turns out to be unappropriate in the case of  $c(p)$  and even unapplicable in the second case ( $(a_n)$  grows too fast). Now we are going to define the generalized Fourier transform as follows.

Let

$$\begin{aligned} p^*(v) &= \inf_{x<0} (p(x) - xv), \text{ if } p \text{ is of concave type;} \\ p^*(v) &= \inf_{x>0} (p(x) - xv), \text{ if } p \text{ is of convex type,} \end{aligned}$$

and we define a function  $r(v)$  by the equality

$$p^*(v) = p(r(v)) - r(v)v.$$

Then for  $a \in c(p)$

$$\tilde{\mathcal{F}}a(z) \stackrel{\text{def}}{=} \sum_{n=r(\log \frac{1}{|z|})}^{+\infty} a_n z^n, \quad |z| < 1;$$

for  $a \in l(p)$

$$\tilde{\mathcal{F}}a(z) \stackrel{\text{def}}{=} \sum_{n=-\infty}^{r(\log \frac{1}{|z|})} a_n z^n, \quad |z| < 1.$$

**Lemma 2.3.** 1) For  $a \in c(p)$

$$\forall c_2 \exists c_1 |\bar{\partial}(\tilde{\mathcal{F}}a)(z)| \leq c_1 e^{-p^*(c_2(1-|z|))}, \quad |z| \sim 1.$$

2) For  $a \in l(p)$

$$\forall c_2 \exists c_1 |\bar{\partial}(\tilde{\mathcal{F}}a)(z)| \leq c_1 e^{-p^*(c_2 \log \frac{1}{|z|})}, \quad |z| \sim 0.$$

Let us remark that for concave-type  $p$  the function  $p^*$  grows near  $v = 0$  and for convex-type  $p$  its  $p^*$  grows towards  $v = \infty$ . Conditions (8) and (9) guarantee that  $\tilde{\mathcal{F}}a$  satisfies either conditions of Theorem 1.1 or conditions of Theorem 1.3. In other words  $\tilde{\mathcal{F}}a$  is AH in the disc near the circle in the first case and AH near the origin in the second case.

Now suppose that

$$a, b \in l(p), \quad (a * b)_n = 0, \quad n < 0, \quad (a * b)_0 = 1. \quad (10)$$

$\tilde{\mathcal{F}}$  is asymptotically multiplicative, which means that

$$|\tilde{\mathcal{F}}(a * b)(z) - \tilde{\mathcal{F}}a(z) \cdot \tilde{\mathcal{F}}b(z)| \rightarrow 0, \quad |z| \rightarrow 0. \quad (11)$$

Denoting  $f \stackrel{\text{def}}{=} \tilde{\mathcal{F}}a$ ,  $g \stackrel{\text{def}}{=} \tilde{\mathcal{F}}b$  and taking into account (10) and (11) we see that (S3) cannot occur neither for  $f$  nor for  $g$ .

Thus (B3) holds and so holomorphic  $f^*$  and  $g^*$  exist such that

$$\frac{1}{2} < |f^*(z)g^*(z)| < 2, \quad |z| \sim 0$$

$$\exists c_1, c_2 : |\log f^*(z)| + |\log g^*(z)| \leq c_1 \sum_{n \geq 0} e^{-c_2 p_n} |z|^{-n}, \quad |z| \sim 0. \quad (13)$$

We need to prove that there exists  $k \in \mathbb{Z}$  such that  $a_m = 0$ ,  $n < k$  and  $b_n = 0$ ,  $n < -k$ . This a reformulation of the desirable standardness of all  $\tau$ -invariant subspaces. Suppose this is not the case. Then  $f$  or  $g$  is not bounded. It is easy to see that  $f^*$  or  $g^*$  is not bounded either. Suppose this for  $f^*$ . Then denoting  $\Phi(z) = \log f^*(z)$  we can assume that  $|\Phi(z)| \geq c|z|^{-1}$  and taking (13) into account we see that

$$cr^{-1} \leq \sum_{n \geq 0} e^{-c_2 p_n} r^{-n}.$$

Now the second condition in (9) allows to show (see [4] e.g.) that

$$p_n \leq nC + \frac{1}{c_2} n \log n$$

which contradicts the first condition in (9).

The proof of Theorem 2.1 is based on the “inner-outer” factorization in  $c(p)$ . It turns out that every  $f \in c(p)$ ,  $f \neq 0$ , can be represented in a form

$$f = Ih,$$

where  $h \in H^2$  and  $I$  is a unimodular function in  $c(p)$ .

### 3. Dynamical System Application

Here we consider the system of differential equations

$$\begin{cases} \dot{x} = \alpha(x, y) \\ \dot{y} = \beta(x, y), \quad (x, y) \in R^2 \end{cases} \quad (14)$$

where  $\alpha, \beta$  are real functions belonging to a very smooth Carleman class  $C\{M_n\}$ . A limit cycle means that there are no other cycles in a neighbourhood. It is proved now by Yu. S. Il'yashenko [8], J. Ecalle and J. Martinet, R. Moussu, J.-P. Ramis [9], [10], [11], that for real analytic  $\alpha, \beta$  in a domain  $\Omega$  there is only a finite number of limit cycles of (14) in any compact part of  $\Omega$ .

The AH technique allows to prove a similar result for quasianalytic vector fields (14), at least if all the singular points of (14) are non-degenerate. To state the result we need some notations.

For a given sequence  $\{M_n\}$  the Carleman class  $C\{M_n\}$  is defined by

$$C\{M_n\} = \{f \in C^\infty : \exists B_f, C_f : |f^{(n)}(x)| \leq B_f C_f^{|n|} M_{|n|}\}.$$

Here  $n = (n_1, n_2)$ ,  $|n| = n_1 + n_2$ . We consider only Carleman classes with regular sequences  $\{M_n\}$ , namely if  $m_n = M_n/n!$  then we assume that

- 1)  $m_n^2 \leq m_{n-1}m_{n+1}$ ,
- 2)  $\sup \left( \frac{m_{n+1}}{m_n} \right)^{1/n} < \infty$ ,
- 3)  $\lim m_n^{1/n} = \infty$ .

Under these assumptions the properties of  $C\{M_n\}$  can be reflected adequately by

$$w(x) \stackrel{\text{def}}{=} \inf_{n \geq 0} x^n \frac{M_n}{n!}.$$

For instance,  $C\{M_n\}$  is quasianalyticique ( $QA$ ) iff

$$\int_0^\infty \log \log \frac{1}{w(x)} dx = \infty \quad (15)$$

Adopting the ideas of Il'yashenko [12] to this new situation it is not hard to prove the following result [13].

**Theorem 3.1.** *Let  $\alpha, \beta$  belong to  $C_\Omega\{M_n\}$  and*

$$w(x) \leq \exp(-\exp \frac{1}{x^{16}}), \quad (16)$$

*then in any compact part of  $\Omega$  there is only a finite number of limit cycles providing that all singular points of (14) are non-degenerate.*

Proof repeats the one in [12] and proceeds by extending the monodromy transform to an AH function in a complex domain. Then Theorem 1.2 is applied. It is very likely that (16) can be weakened to (15).

#### 4. Weighted Polynomial Approximation

Much interesting analysis has resulted from attempts to understand the structure of  $P^2(\mu)$ , the closure in  $L^2(\mu)$  of the set  $\mathcal{P}_A$  of all analytic polynomials, where  $\mu$  is a positive finite Borel measure with compact support in the complex plane  $\mathbb{C}$ . The greatest achievement in this field due to J. Thomson [14] asserts that  $P^2(\mu) \neq L^2(\mu)$  if and only if there is a point  $c \in \mathbb{C}$  such that

$$|p(c)| \leq k \|p\|_{L^2(\mu)}, \quad \forall p \in \mathcal{P}_A.$$

Such points are called points of bounded point evaluation. In other words Thomson has solved a problem of Mergelyan-Brennan. He even has achieved more: a description of  $P^2(\mu)$  in terms of points of bounded evaluation.

But unfortunately this does not help much when one is interested in a description of  $P^2(\mu)$  in terms of  $\mu$  itself. Here are two examples. The first is a so-called splitting problem. In this problem

$$d\mu = w(1-r)rdrd\theta + h(\theta)dm(\theta) = \mu_D + \mu_T$$

is a sum of a radially symmetric measure  $\mu_{\mathbb{D}}$  in the disc  $\mathbb{D}$  and a measure  $\mu_T$  on  $T = \partial\mathbb{D}$ . What are the conditions on the pair  $(w, h)$  necessary and sufficient for the splitting:

$$P^2(\mu) = P^2(\mu_{\mathbb{D}}) \oplus L^2(\mu_T)?$$

This problem was solved in [15], [16] and the next result is nothing more than another form of Theorem 1.1.

**Theorem 4.1.** 1) Suppose that (R1) holds. Then the following conditions are sufficient for the splitting

$$\int_T \log h dm = -\infty, \quad \int_0 \log \log \frac{1}{w(x)} dx = \infty. \quad (17)$$

2) Conditions (17) are also necessary if there exists an arc  $I$ ,  $I \in T$ , such that

$$\int_I h^{-1} dm < \infty.$$

In our second example we deal with a problem of weighted polynomial approximation in an arbitrary simply connected domain. Let  $\Omega$  be such a domain,  $u$  is a positive continuous function in  $\Omega$  vanishing towards the boundary. It is always the case that

$$P^2(udm_2) \subset L_a^2(udm_2) \stackrel{\text{def}}{=} L^2(udm_2) \cap \text{Hol}(\Omega).$$

Keldysh was the first who investigated when the equality

$$P^2(udm_2) = L_a^2(udm_2) \quad (18)$$

holds. He considered  $u(z) = u(\text{dist}(z, \partial\Omega))$ . This is not convenient as the conditions on  $U$  depend heavily on the smoothness of  $\partial\Omega$  and vary from

$$u(x) \leq \exp \left( -\exp \frac{1}{x^2} \right)$$

to

$$u(x) \leq \exp \left( -\exp \frac{1}{x} \right).$$

J. Brennan considered the case  $u(z) = w(G(z))$  where  $G$  is the Green function of  $\Omega$ . The result of Brennan [17] is the following.

**Theorem 4.2.** 1) Let  $\Omega$  be an arbitrary simply connected domain,  $G$  be its Green function,

$$u(z) = w(G(z)),$$

where  $w$  satisfies (R1). Then the condition

$$\int_0 \log \log \frac{1}{w(x)} dx = \infty$$

is sufficient for (18).

2) If  $\partial\Omega$  is smooth enough, then this condition is also necessary.

The proof is based on Theorem 1.1. Under some auxiliary conditions on  $\Omega$  this theorem was proved in [18].

## References

- [1] Domar, Y: A solution of the translation-invariant subspace problem for weighted  $L^p$  on  $\mathbb{R}$ ,  $\mathbb{R}_+$  or  $\mathbb{Z}$ . Lecture Notes in Mathematics, vol. 975. Springer, Berlin Heidelberg New York 1983, pp. 214–226
- [2] Domar, Y: Extension of the Titchmarsh convolution theorem with applications in the theory of invariant subspaces. Proc. London Math. Soc. 1983, no. 46, pp. 288–300
- [3] Domar, Y: Translation invariant subspaces of weighted  $L^p$  and  $l^p$  spaces. Math. Scand. **49** (1981) 133–144
- [4] Borichev, A.A: A Titchmarsh-type convolution theorem on the group  $\mathbb{Z}$ . Arkiv for Matematik **27** (1989) 179–187
- [5] Borichev, A.A: Generalized Fourier transform, the Titchmarsh theorem and almost analytic function. Algebra and Analysis, no. 4 (1989) 17–53 (in Russian; translation in Leningrad J. Math.)
- [6] Sarason, D: The  $H^p$  spaces of an annulus. Mem. Amer. Math. Soc. **56** (1965) 1–78
- [7] Hitt, D: Invariant subspaces of  $H^2$  of an annulus. Pac. J. Math. **134** (1988) 101–133
- [8] Il'yashenko, Yu. S: Finiteness theorems for limit cycles. Usp. Mat. Nauk **45**, no. 2 (1990) 143–200 (in Russian; translation in Russ. Math. Surv.)
- [9] Ecalle, J., Martinet J., Moussu R., Ramis, J.-P: Non-accumulation des cycles-limites (I). C.R. Acad. Sci. Paris **304**, no. 13 (1987) 375–377
- [10] Ecalle, J., Martinet, J., Moussu, Rl, Ramis, J.-P: Non-accumulation des cycles-limites (II). C.R. Acad. Sci. Paris **304**, no. 14 (1987) 431–434
- [11] Ecalle, J: Finitude des cycles-limites et accelero-sommation de l'application de retour. Pré publications, Université de Paris-Sud, ORSAY, 1990, no. 90–36, pp. 1–86
- [12] Il'ashenko, Yu.S: Limit cycles of polynomial vector fields with non-degenerate singular points on the real plane. Funk. Anal. Priloz. **18**, no. 3 (1984) 32–42 (in Russian; translation in Funct. Anal. Appl., January 1985, pp. 199–209)
- [13] Volberg, A.L: Un théorème de Dulac-Ecalle-Il'yashenko-Martinet-Moussu-Ramis etendu aux fonction quasianalytique. Publ. Math. d'ORSAY, Séminaire d'Analyse Harmonique, Année 1989/90, pp. 152–171
- [14] Thomson, J: Approximation in the mean by polynomials. To appear in Ann. Math.
- [15] Volberg A.L: The logarithm of an almost analytic function is summable. Sov. Math. Dokl. **26** (1982) 238–243
- [16] Volberg A.L., Jöricke, B: Summability of the logarithm of an almost analytic function and a generalizations of the Levinson-Cartwright theorem. Matem. Sbornik **130** (1986) 335–348 (in Russian; translation in Math. USSR Sbornik **58** (1987) 337–349)
- [17] Brennan, J: Weighted polynomial approximation and quasianalyticity for general sets. Preprint, University of Kentucky, 1990
- [18] Brennan, J: Weight polynomail approximation, quasianalyticity and analytic continuation. J. Reine Angew. Math. **357** (1985) 23–50



# Cyclic Cohomology and $K$ -Homology

*Joachim Cuntz*

Mathematisches Institut, Universität Heidelberg, Im Neuenheimer Feld 288  
W-6900 Heidelberg, Fed. Rep. of Germany

Cyclic cohomology, as introduced by A. Connes (and independently by Tsygan), can be viewed as an algebraic version of  $K$ -homology in the sense of [BDF], [Kas]. In [Cu1], the author gave a description of  $K$ -homology for a complex algebra  $A$  using the free product algebra  $QA = A * A$  and its ideal  $qA$  obtained as the kernel of the multiplication map  $QA \rightarrow A$ . The even  $K$ -homology group  $K^0 A$  can be defined as the set  $[qA, K]$  of homotopy classes of homomorphisms from  $qA$  into the elementary algebra  $K$  of compact operators (the completion of the algebra of complex matrices of arbitrary size). This set is an abelian group in a natural way. To describe the odd  $K$ -homology group  $K^1 A$  one can use the universal extension  $0 \rightarrow qA \rightarrow RA \rightarrow A \rightarrow 0$ , where  $RA$  is the free non-unital tensor algebra over  $A$ , described below in Section 1. In contrast with the situation for  $qA$ , the algebra  $qA$  admits several different  $C^*$ -algebra completions if  $A$  is a  $C^*$ -algebra. Two natural completions  $\overline{qA}^\ell$ ,  $\overline{qA}^{cp}$  can be defined by the condition that the maps  $\overline{RA}^\ell \rightarrow A$ ,  $\overline{RA}^{cp} \rightarrow A$  admit a linear, resp. completely positive splitting of norm 1. With these completions,  $[\overline{qA}^\ell, K]$  is a version of the Brown-Douglas-Fillmore group  $\text{Ext}_A$ , while  $[\overline{qA}^{cp}, K]$  is the Kasparov group  $K^1 A$ . The set  $[\overline{qA}^\ell, K]$  is only a semigroup in general. More generally, these algebras can be used to define the Kasparov groups  $KK_*(A, B)$  as  $KK_0(A, B) = [qA, K \otimes B]$ ,  $KK_1(A, B) = [\overline{qA}^{cp}, K \otimes B]$ . The algebra  $\overline{qA}^{cp}$  is in fact isomorphic to the subalgebra of even elements in the natural  $C^*$ -completion of  $qA$  (it is a hereditary subalgebra of the algebra  $\varepsilon A$  used by Zekri in [Zek]). The basic result of  $KK$ -theory – the existence of the Kasparov product – is in this setting the homotopy equivalence between  $K \otimes qA$  and  $K \otimes q(qA)$  [Cu1].

We now come to the connection with cyclic cohomology. It was shown in [Co-Cu] that cyclic cocycles, viewed as cocycles in the bicomplex of Tsygan-Loday-Quillen or in the  $B - b$  bicomplex of Connes, correspond exactly to traces or supertraces on a power  $(qA)^n$  of the ideal  $qA$  or to supertraces on  $QA$  that vanish on  $(qA)^n$ . In [Cu-Qu], among other things, it is shown that such a supertrace corresponds to a cyclic coboundary if and only if it is of the form  $T' \circ d$  where  $d : QA \rightarrow \Omega^1 QA$  is the universal derivation into the bimodule of abstract 1-forms over  $QA$  and  $T'$  is a supertrace on this bimodule. In fact, composition with  $d$  basically is the coboundary operator. Thus a cyclic cohomology class is roughly speaking a “homotopy class” of supertraces on  $QA$  or  $(qA)^n$ .

For this correspondence between supertraces and cocycles-coboundaries to be natural one has to introduce certain constants into the basic bicomplex (see

3.A and 3.B below). These constants become important when constructing cyclic cycles in the bicomplex. By what has been said above, such a cycle determines a “closed” element of  $QA$ , i.e. an element  $x$ , that varies under homotopies only by sums of commutators or, formally, for which  $dx$  is a sum of commutators in  $\Omega^1 QA$ . The natural cycles associated with idempotents or invertible elements in  $A$  are represented by certain power series in  $QA$  (whose coefficients are computed using the constants in the bicomplex) that have quite natural interpretations.

The exposition below is largely based on joint work with D. Quillen [Cu-Qu] and also on [Co-Cu, Cu2]. Needless to say, much of the material refines and continues ideas of A. Connes [Co1, Co2].

## 1. Universal Algebras Associated with an Algebra $A$

Let  $A, B$  be algebras and  $E$  a bimodule over  $B$ . We will see below that multilinear maps  $\omega : A^n \rightarrow E$  satisfying

$$\begin{aligned} p(x_0)\omega(x_1, \dots, x_n) - \omega(x_0x_1, \dots, x_n) + \omega(x_0, x_1x_2, \dots, x_n) \\ - \dots + (-1)^n\omega(x_0, \dots, x_{n-1}x_n) + (-1)^{n+1}\omega(x_0, \dots, x_{n-1})p(x_n) = 0 \end{aligned} \quad (1.A)$$

where  $p : A \rightarrow B$  is a linear map, together with a similar multilinear map  $\omega'$  into the dual of  $E$  lead to cyclic cocycles on  $A$ . Let us see first how such maps may arise:

**0.** An example for  $n = 0$  is given by a central element  $z \in B$  and a linear map  $p : A \rightarrow B$ . One has  $p(x)z - zp(x) = 0$ .

**1.** In the one-variable case the main example is given by a quasihomomorphism, that is a pair of homomorphisms  $\alpha, \bar{\alpha} : A \rightarrow B'$  into some algebra  $B'$  that contains  $B$  as an ideal such that  $\alpha(x) - \bar{\alpha}(x) \in B$ ,  $\forall x \in A$ . Putting  $p(x) = \frac{1}{2}(\alpha(x) + \bar{\alpha}(x))$ ,  $q(x) = \frac{1}{2}(\alpha(x) - \bar{\alpha}(x))$  we have  $q(xy) = p(x)q(y) + q(x)p(y)$ . More generally, for  $\alpha, \bar{\alpha}$  we may take two linear maps whose curvatures (see 2. below) coincide.

**2.** For the canonical example with 2 variables let  $\varphi : A \rightarrow B$  be a linear map between two algebras. Let  $\omega(x, y) = \varphi(xy) - \varphi(x)\varphi(y)$  be its “curvature”, cf. [Qu]. Then  $\varphi(x)\omega(y, z) - \omega(xy, z) + \omega(x, yz) - \omega(x, y)\varphi(z) = 0$ .

Finally, if  $E = B$  and  $\omega_1, \omega_2$  are multilinear maps of  $n_1$  resp.  $n_2$  variables satisfying (1.A) then the product  $\omega_1\omega_2$  is a function of  $n_1 + n_2$  variables satisfying (1.A).

We now introduce, given an algebra  $A$ , algebras  $RA, QA$  with a linear map  $p : A \rightarrow RA \subset QA$  in which these examples are universally realized. The algebra  $RA$  is simply the free tensor algebra over  $A$ ,  $RA \cong \bigoplus_{n \geq 1} A^{\otimes n}$  or, in other words the universal algebra generated by symbols  $p(x)$ ,  $x \in A$  which are linear in  $x$  with no further relations ( $p(x_1) \dots p(x_n)$  corresponding to  $x_1 \otimes x_2 \otimes \dots \otimes x_n$  in the tensor algebra). Let

$$\omega(x, y) = p(xy) - p(x)p(y)$$

denote the curvature of the universal linear map  $p$  and  $\varrho A$  the ideal in  $RA$  generated by  $\omega(x, y)$ ,  $x, y \in A$ . There is a short exact sequence

$$0 \rightarrow \varrho A \rightarrow RA \rightarrow A \rightarrow 0.$$

This is a (uni)versal extension of  $A$ : if  $0 \rightarrow J \rightarrow B \rightarrow A \rightarrow 0$  is any extension with a linear splitting  $\varphi : A \rightarrow B$  then the map  $p(x) \mapsto \varphi(x)$  induces a homomorphism  $\tilde{\varphi} : RA \rightarrow B$  that maps  $qA$  into  $J$ . Also, the algebra  $RA$  is contractible via the homotopy  $\tilde{\varphi}_t$ ,  $\varphi_t : p(x) \mapsto tp(x)$  so that  $RA$  may be thought of as a non-commutative cone over  $A$ , the ideal  $qA$  corresponding to the suspension.

The algebra  $QA$  now is defined as the universal algebra generated by symbols  $p(x)$ ,  $q(x)$ ,  $x \in A$  which are linear in  $x$  and satisfy

$$\begin{aligned} q(xy) &= p(x)q(y) + q(x)p(y) \\ q(x)q(y) &= p(xy) - p(x)p(y). \end{aligned}$$

There are two natural homomorphisms  $\iota, \bar{\iota} : A \rightarrow QA$  mapping  $x$  to  $p(x) \pm q(x)$  and a natural  $\mathbb{Z}/2$ -grading of  $QA$  with grading automorphism  $-$  defined by  $(p(x) + q(y)) = p(x) - q(y)$ . In fact,  $QA$  is isomorphic to the free product  $A * A$  with the grading automorphism that exchanges the two copies of  $A$ .

**Proposition 1.1.** *The natural map  $RA \rightarrow QA$  is injective and induces an isomorphism between  $RA$  and the algebra  $QA_+$  of even elements in  $QA$ .*

In  $QA$ , let  $qA$  be the ideal generated by  $q(x)$ ,  $x \in A$ . There is a short exact sequence

$$0 \rightarrow qA \rightarrow QA \rightarrow A \rightarrow 0$$

with two different splittings  $A \rightarrow QA$ . The map in 1.1 maps  $qA$  onto  $((qA)^2)_+$ . There is a simple argument showing that  $QA$  is homotopy equivalent in  $2 \times 2$ -matrices to the direct sum  $A \oplus A$  (cf. [Cu1]). If we filter  $QA$  by the powers of the ideal  $qA$ , then the associated graded algebra is simply the algebra  $\Omega A$  of abstract differential forms over  $A$ , i. e. the universal algebra generated by symbols  $\varphi(x)$  and  $d(x)$  where  $\varphi, d$  are linear in  $x$  and satisfy the relations

$$\begin{aligned} \varphi(xy) &= \varphi(x)\varphi(y) \\ d(xy) &= \varphi(x)d(y) + d(x)\varphi(y). \end{aligned}$$

Just as  $qA$  is the universal ideal in an extension of  $A$ ,  $qA$  is universal for quasihomomorphisms as in 1. above: given a pair of homomorphisms  $\alpha, \bar{\alpha} : A \rightarrow B'$  where  $B'$  contains  $B$  as an ideal, we get a homomorphism  $QA \rightarrow B'$  mapping  $p(x)$ ,  $q(x)$  to  $\frac{1}{2}(\alpha(x) + \bar{\alpha}(x))$ ,  $\frac{1}{2}(\alpha(x) - \bar{\alpha}(x))$ . This homomorphism sends  $qA$  into the ideal  $B$ . Quasihomomorphisms arise for instance from even Fredholm or Kasparov modules.

## 2. Transport of K-Theory and Cyclic Cohomology Classes by Certain Linear Maps

K-Theory elements may be transported not only by homomorphisms but also by certain linear (or even multilinear) maps between algebras. Let, for instance,  $\varphi : A \rightarrow B$  be a linear map and  $\omega(x, y) = \varphi(x, y) - \varphi(x)\varphi(y)$ ,  $x, y \in A$  its curvature. If  $\varphi$  is a  $*$ -respecting map between  $C^*$ -algebras,  $e = e^2 = e^*$  is a projection, and  $\|\omega(e, e)\| < \frac{1}{2}$  then  $\varphi(e)$  is close to a projection in  $B$  that can be obtained from  $\varphi(e)$  by functional calculus, cf. also [CGM]. Similarly, if  $A, B$  are

arbitrary algebras and  $\omega$  is small algebraically e.g. nilpotent, then  $\varphi$  can be used to transport algebraic  $K$ -theory classes.

As a second example, consider a bounded linear map  $\varphi : A \rightarrow B$  of two  $C^*$ -algebras for which  $\varphi(x)\omega(y, z)$ ,  $\omega(y, z)$  for  $x, y, z \in A$  lie in some closed subalgebra  $B_0 \subset B$ . To show that  $\varphi$  induces a map  $K_*A \rightarrow K_{*+1}(B_0)$  we use the algebras  $\varrho A$ ,  $RA$  completed with respect to the maximal  $C^*$ -norm for which  $p$  is of norm  $\leq \|\varphi\|$ . We obtain an exact sequence of  $C^*$ -algebras

$$0 \rightarrow \varrho A \rightarrow RA \rightarrow A \rightarrow 0$$

where  $RA$  is contractible. The long exact sequence of  $K$ -theory yields an isomorphism

$$K_*(A) \xrightarrow{\cong} K_{*+1}(\varrho A).$$

Since  $\varphi$  induces a homomorphism  $\tilde{\varphi} : \varrho A \rightarrow B_0$ , we obtain the desired map

$$K_*(A) \cong K_{*+1}(\varrho A) \rightarrow K_{*+1}(B_0).$$

If  $\varphi$  is completely positive then it induces even a  $KK$ -element since then, for an appropriate completion,  $\varrho A$  becomes a hereditary subalgebra of the  $C^*$ -algebra completion of  $\varepsilon A = \varrho A \times \mathbb{Z}/2$  on which  $\tilde{\varphi}$  is defined.

Consider further a quasihomomorphism, i.e. a pair  $\alpha, \bar{\alpha} : A \rightarrow B$  such that  $\alpha(x) - \bar{\alpha}(x) \in J$ ,  $x \in A$  for some ideal  $J$  of  $B$ . We assume that  $A, B, J$  are  $C^*$ -algebras, even though part of the discussion goes through more generally. This quasihomomorphism induces a homomorphism  $QA \rightarrow B$  restricting to  $\bar{\alpha} : \varrho A \rightarrow J$ . By split exactness of  $K_*$ ,  $K_*(QA) \cong K_*(\varrho A) \oplus K_*(A)$  and the map  $K_*(i) - K_*(\bar{i}) : K_*(A) \rightarrow K_*(QA)$ , where  $i, \bar{i}$  are as in Section 1, sends  $K_*(A)$  into the first factor  $K_*(\varrho A)$ . Composing with  $K_*(\bar{\alpha})$  yields a map  $K_*(A) \rightarrow K_*(J)$ . We note that Kasparov or Fredholm modules usually give rise to quasihomomorphisms, thus can be used to transport  $K$ -theory classes.

We will see below that cyclic cocycles can be described as traces on  $RA$  or  $QA$ . Since quasihomomorphisms or linear maps from  $A$  into some algebra  $B$  induce homomorphisms  $QA, RA \rightarrow B$  it will follow that traces on  $B$  can be transported under such maps to cyclic cocycles on  $A$ . The pairing between a trace and a transported  $K$ -theory class corresponds to the pairing between the  $K$ -theory class and the transported trace.

### 3. Traces on $RA$ , Supertraces on $QA$ and Their Homotopy Classes

Let  $QA_+$ ,  $QA_-$  denote the sets of even and odd elements with respect to the  $\mathbb{Z}/2$ -grading  $QA$ , respectively. We have seen above that  $QA_+ \cong RA$  as an algebra.  $QA_+$  and  $QA_-$  are  $RA$ -bimodules. Recall that a supertrace on a  $\mathbb{Z}/2$ -graded algebra or bimodule is a linear functional  $T$  satisfying  $T(x\alpha) = (-1)^{\deg(x)\deg(\alpha)} T(\alpha x)$ . Every even (resp. odd) supertrace on  $QA$  gives a bimodule trace  $T$  with  $T(\alpha\omega) = T(\omega\alpha)$  for  $\alpha \in RA$ ,  $\omega \in QA_+$  (resp.  $\omega \in QA_-$ ). Conversely, a bimodule trace  $T$  comes from a supertrace if and only if  $T(qx_0 \dots qx_n) = T(qx_n qx_0 \dots qx_{n-1})$  for  $x_0, \dots, x_n \in A$ .

A linear functional on  $((qA)^{2m})_+ = (\varrho A)^m$  (an important special case is of course  $(\varrho A)^0 = RA$ ) is given by its components  $T^n$ ,  $n \geq 2m$ , defined by

$$T^n(x_0, \dots, x_k) = \begin{cases} T(qx_0qx_1\dots qx_n) & n \text{ odd} \\ T(px_0qx_1\dots qx_n) & n \text{ even.} \end{cases}$$

**Proposition 3.1.**  *$T$  represents a (bimodule)-trace if and only if for all even  $n \geq 2m$*

$$b T^n = (1 + \lambda) T^{n+1}$$

$$b' T^{n-1} = (1 - \lambda) T^n$$

(here  $\lambda, b, b'$  are defined as in [Co1]). Every bimodule trace can be modified canonically so as to extend to an even supertrace  $\tilde{T}$  on  $(qA)^{2m}$ .

To prove Proposition 3.1, consider the map  $\delta : QA_{\pm} \rightarrow QA_{\mp}$  defined by  $\delta(\omega q(x)) = [p(x), \omega]$ ,  $\omega \in QA$ . The image of  $\delta$  consists of all commutators in the bimodules  $QA_+$ ,  $QA_-$  and one checks easily for  $\delta : QA_- \rightarrow QA_+$

$$\begin{aligned} (T \circ \delta)^{2k+1} &= -b T^{2k} + (1 + \lambda) T^{2k+1} \\ (T \circ \delta)^{2k} &= -b' T^{2k-1} - (1 - \lambda) T^{2k}. \end{aligned}$$

Since  $(1 - \lambda^2)T^k = (1 - \lambda)bT^{k-1} = b'(1 - \lambda)T^{k-1} = b'b'T^{k-2}$ , the relations of Proposition 3.1 imply that  $(1 - \lambda^2)T^k = 0$  for odd  $k \geq 2n + 1$ , so that  $(1 + \lambda)T^k = \frac{2}{k+1}N T^k$  where  $N = 1 + \lambda + \dots + \lambda^k$ . Thus  $T$  is a trace if and only if  $(1 - \lambda^2)T^k = 0$  for such  $k$  and  $T$  is a cocycle in the total complex of the bicomplex

$$\begin{array}{ccc} C^{2n+1} & \xrightarrow{\frac{-2}{2n+2}N} & \\ \uparrow \frac{2n+2}{2}b' & & \uparrow b \\ \xrightarrow{N} & C^{2n} & \xrightarrow{1-\lambda} \\ \uparrow -nb & & \uparrow -b' \\ \xrightarrow{-n(1-\lambda)} & C^{2n-1} & \end{array} \quad (3.A)$$

It is obvious that every cocycle  $T$  in this bicomplex is cohomologous to one for which  $(1 - \lambda)T^k = 0$  for all odd  $k$ . This procedure to obtain cyclic cocycles admits the following generalization:

Let  $A$  be an algebra and  $E, E'$  two  $RA$ -bimodules which are in duality, i.e. there is a bilinear pairing  $(x, x') \mapsto \langle x|x' \rangle \in \mathbb{C}$ ,  $x \in E$ ,  $x' \in E'$  such that  $\langle \alpha x|x' \rangle = \langle x|x'\alpha \rangle$  and  $\langle x\alpha|x' \rangle = \langle x|\alpha x' \rangle$  for  $\alpha \in RA$ . Denote as above by  $p : A \rightarrow RA$  the universal linear map and by  $\omega$  its curvature. Suppose further that we are given two multilinear maps  $c : A^r \rightarrow E$ ,  $c' : A^s \rightarrow E'$  satisfying

$$\begin{aligned} p(x_0)c(x_1, \dots, x_r) - c(x_0x_1, \dots, x_r) + \dots \\ + (-1)^r c(x_0, x_1, \dots, x_{r-1}x_r) + (-1)^{r+1} c(x_0, \dots, x_{r-1})p(x_r) = 0 \end{aligned}$$

and similarly for  $c'$  with  $r$  replaced by  $s$ . Then the following family  $\{T^n\}_{n \geq r+s}$  of multilinear functions defines a cocycle in the bicomplex 3.A resp. 3.B:

$$T^n(x_0, \dots, x_n) = \begin{cases} \langle \omega(x_0, x_1) \dots \omega(x_{k-1}, x_k) c(x_{k+1}, \dots, x_{k+r}) | c'(x_{k+r+1}, \dots, x_n) \rangle & k \text{ odd} \\ \langle px_0 \omega(x_1, x_2) \dots \omega(x_{k-1}, x_k) c(x_{k+1}, \dots, x_{k+r}) | c'(x_{k+r+1}, \dots, x_n) \rangle & k \text{ even} \end{cases}$$

where  $k = n - r - s$ . In fact, by the same computation as above we find that  $b' T^n = (1 + \lambda) T^{n+1}$  for  $k$  even, and  $b' T^n = (1 - \lambda) T^{n+1}$  for  $k$  odd.

Applying this with  $c(x_1, \dots, x_{2n+1}) = q(x_1)q(x_2)\dots q(x_{2n+1})$  we see that for bimodule traces on  $((QA)^{2n+1})_-$ , Proposition 3.1 holds with odd and even interchanged. A linear functional on  $((QA)^{2n+1})_-$  is a trace if and only if  $(1 - \lambda^2) T^k = 0$  for even  $k \geq 2n + 1$  and if  $T$  is a cocycle in the bicomplex

$$\begin{array}{ccccc}
C^{2n+2} & \xrightarrow{\frac{-2}{2n+3}N} & & & \\
\uparrow \frac{2n+3}{2}b' & & \uparrow b & & \\
& \longrightarrow & C^{2n+1} & \xrightarrow{1-\lambda} & \\
& & \uparrow \frac{-(2n+1)}{2}b & & \uparrow -b' \\
& & \xrightarrow{\frac{-(2n+1)}{2}(1-\lambda)} & & C^{2n}
\end{array} \tag{3.B}$$

The question of which traces correspond to coboundaries now admits a very natural answer which is a basis for many different homotopy invariance results for cyclic cohomology. We form the algebra  $\Omega Q A$  of abstract differential forms over  $Q A$ . An important fact which is easily checked is that the map  $d(p(x)), d(q(x)) \mapsto p(dx), q(dx)$  induces an isomorphism (of bigraded algebras) between  $\Omega Q A$  and  $Q \Omega A$ . Let  $\Omega^1 Q A$  be the  $Q A$ -bimodule of 1-forms over  $Q A$  and  $d : Q A \rightarrow \Omega^1 Q A$  the universal derivation. If  $T'$  is a (bimodule) supertrace on  $\Omega^1 Q A$  then  $T' \circ d$  is a supertrace on  $Q A$ .

We now determine the supertraces  $T'$  on  $\Omega^1 Q A$  and show that the cocycles corresponding to  $T' \circ d$  are coboundaries. We restrict to the case of even supertraces. The odd case is quite similar but there are some slight complications in dimension 0 and 1 (an odd supertrace  $T$  corresponds to a cocycle in the bicomplex 3.B only if  $T(q(x)) = 0, \forall x$ ) which we want to avoid.

**Proposition 3.2.** *There is a bijection between even supertraces  $T'$  on  $\Omega^1 Q A$  and cochains  $\alpha = \{\alpha^k\}_{k \geq 0}$  given by*

$$\begin{aligned}
\alpha^{2n}(x_0, \dots, x_{2n}) &= T'(dx_0qx_1\dots qx_{2n}) \\
\alpha^{2n+1}(x_0, \dots, x_{2n+1}) &= T'(x_0dx_1qx_2\dots qx_{2n+1})
\end{aligned}$$

Let  $\partial^{(1)}, \partial^{(2)}$  denote the two coboundary operators in the bicomplex 3.A, i.e.

$$\begin{aligned}
\partial_{2n}^{(1)} &= N + (n+1)b', \quad \partial_{2n-1}^{(1)} = -n(b + (1-\lambda)) \\
\partial_{2n}^{(2)} &= b + 1 - \lambda, \quad \partial_{2n-1}^{(2)} = -(b' + \frac{1}{n}N)
\end{aligned}$$

One has  $\partial^{(2)}\partial^{(1)} = \partial^{(1)}\partial^{(2)} = 0$  and we already have seen in 3.1 that a cochain  $T$  represents an even supertrace on  $Q A$  if and only if  $\partial^{(2)} T = 0$  and  $T$  is normalized, i.e.  $(1 - \lambda) T^{2n+1} = 0, \forall n$ .

Every cochain  $\alpha$  can be modified canonically to a cochain  $P\alpha$  such that  $\partial^{(1)}(P\alpha)$  is normalized. One has

$$(P\alpha)^{2n} = \alpha^{2n}$$

$$(P\alpha)^{2n-1} = \frac{1}{2n} N \alpha^{2n-1} + (1-\lambda)^{-1} \left(1 - \frac{1}{2n} N\right) b' \alpha^{2n-2}$$

If  $\partial^{(1)}\alpha$  is already normalized then  $P\alpha = \alpha$  so that  $P^2 = P$ . We can now formulate the main result of this section.

**Theorem 3.3** [Cu-Qu]. *Let  $T'$  be an even supertrace on  $\Omega^1 QA$  corresponding to the cochain  $\alpha$ . Then the cocycle corresponding to  $T' \circ d$  is given by  $\partial^{(1)}(P\alpha)$ .*

The analogous result holds for even supertraces on  $(qA)^{2n}$  and for even supertraces on the  $(qA)^{2n}$ -bimodule generated in  $\Omega^1 QA$  by  $p(\Omega^1 A)$ . Theorem 3.3 contains of course as a very special case the fact that cyclic cohomology is invariant under derivations of  $A$  [Go] and has many applications. We list here a few examples:

1. The cyclic cohomology classes associated with an  $n$ -summable Fredholm module or with a  $\Theta$ -summable Fredholm module, [Co1], [Co2] are (in the first case upon application of  $S$ ) invariant under homotopies of the homomorphisms and of the operators appearing in the module. In fact, these classes are obtained from homomorphisms of  $(qA)^n$  or of  $QA$  into an algebra with a trace.

2. A normalized cyclic cocycle is a coboundary if and only if it comes from a supertrace of the form  $T' \circ d$ ,  $T'$  a trace on  $\Omega^1 QA$ . Thus elements of cyclic cohomology can be described as homotopy classes of supertraces on  $QA$  or on  $(qA)^n$ . This applies to all different versions of cyclic cohomology (ordinary, periodic, entire). For instance, the periodic cyclic cohomology  $HC_{\text{per}}^{\text{even/odd}}$  is given as

$\{\text{even/odd supertraces on } \overline{QA}\}/\{T' \circ d \mid T' \text{ an even/odd supertrace on } \Omega^1 \overline{QA}\}$   
where  $\overline{QA} = \varprojlim_n QA/(qA)^n$  is the algebraic completion of  $QA$ .

3. Let  $\varphi : A \rightarrow B$  be a linear map with nilpotent curvature  $(\omega(x_0, x_1)\omega(x_2, x_3) \dots \omega(x_{2k}, x_{2k+1})) = 0$  for  $k \geq n$ ) and  $\tau$  a trace on  $B$ . We obtain a homomorphism  $\tilde{\varphi} : RA \rightarrow B$  and a trace  $T = \tau \circ \tilde{\varphi}$  on  $RA/(qA)^{n+1}$  and thus an element  $[T]$  of  $HC^{2n}(A)$ . Let  $\varphi_t : A \rightarrow B$  be a differentiable homotopy of such maps and  $\dot{\varphi}_t$  its derivative. Consider the linear map  $\Psi_t : \Omega A \rightarrow B$  sending  $x_0 dx_1 \dots dx_n$  to  $\varphi_t(x_0)\varphi_t(x_1) \dots \varphi_t(x_n)$  and the corresponding homomorphism  $\tilde{\Psi}_t : R\Omega A \rightarrow B$ . We obtain a trace  $T'_t$  on  $R\Omega A$  which we can modify to an even supertrace still denoted by  $T'_t$  on  $R\Omega A \cong \Omega QA$  (see 3.1).  $T'_t \circ d$  gives an even supertrace on  $QA/(qA)^{2n+2}$  inducing a coboundary. This shows that the classes  $[T_t]$  induced by  $\varphi_t$  remain constant.

#### 4. Characteristic Classes as Elements of $RA$ and $QA$

The idea of interpreting characteristic classes associated with idempotents or invertible elements as differential forms, elements of  $\Omega A$  [Ka] or even as elements of  $QA$  [Co2] is of course not new. The description of these classes that we are now going to give is however at least partly new and will clarify their meaning.

Let  $e$  be an idempotent in  $A$ . With  $e$  we want to associate odd elements  $ch_{2n}(e) \in (qA)^{2n+1}$  such that  $d(ch_{2n}(e))$  is a sum of supercommutators in  $\Omega^1 QA$  and

such that  $ch_{2n+2k}(e)$  differs from  $ch_{2n}(e)$  by a sum of supercommutators. Second, we associate an element  $Ch(e)$  (an even element of  $\overline{QA}$ ) of the algebraic completion  $RA = \lim_{\leftarrow} RA/(\varrho A)^n$  such that  $d(Ch(e))$  is a sum of supercommutators. The elements  $ch_{2n+2k}(e)$  represent cyclic cycles in  $HC_{2n}^-(A)$  which may be paired with odd traces  $T$  on  $(qA)^{2n}$ . The identities  $T(ch_{2n+2k}(e)) = T(ch_{2n}(e))$ ,  $(T' \circ d)(ch_{2n}(e)) = 0$  mean that we have a pairing with  $\langle [T], ch_{2n}(e) \rangle = \langle S^k[T], ch_{2n-2k}(e) \rangle$ . The element  $Ch(e)$  represents a cycle for  $HC_{\text{ev}}^{\text{per}}(A)$  (or for the entire cyclic homology  $HC_{\text{ev}}(A)$  in the sense of [Co2] or for other versions of cyclic cohomology for locally convex algebras – we only need  $Ch(e)$  to converge in the completion of  $RA$ ) and the pairing with a periodic cocycle represented by a trace  $T$  on  $RA$  is given by  $T(Ch(e))$ . By functoriality, to construct  $ch_k(e)$ ,  $Ch(e)$ , we may assume that  $A = \mathbb{C}e \cong \mathbb{C}$ .

The first part is easily achieved taking  $ch_{2n}(e) = (2qe)^{2n+1}$ . These are elements of  $(qA)^{2n+1}$  differing from the class in dimension 0, given by  $2q(e)$ , only by sums of supercommutators by 3.1 (the difference is a boundary in the bicomplex dual to 3.A).

The element  $Ch(e)$  is given by a series of the form

$$Ch(e) = \sum_{k \geq 0} c_{2k} p(e) q(e)^{2k} + \sum_{k \geq 1} c_{2k-1} q(e)^{2k}.$$

The fact that  $\partial(\sum c_j e^{\otimes^{(j+1)}}) = 0$  in the bicomplex dual to (3.A) leads, given  $c_0 = 1$ , immediately to

$$c_{2k} = -c_{2k-1} = \frac{(2k)!}{(k!)^2}$$

whence, formally,

$$Ch(e) = \frac{p(e)}{\sqrt{1 - 4(qe)^2}} - \frac{1}{2} \left( \frac{1}{\sqrt{1 - 4(qe)^2}} - 1 \right)$$

which gives, using  $(qe)^2 = p(e) - p(e)^2$

$$Ch(e) = \frac{1}{2} \left( 1 - \frac{1 - 2p(e)}{\sqrt{(1 - 2p(e))^2}} \right).$$

This slightly surprising formula has a very natural interpretation: If we replace  $p(e)$  by the complex variable  $z$  (or by an element of a Banach algebra) then  $Ch(z)^2 = Ch(z)$  and  $Ch(z)$  is the characteristic function of the half-plane  $\text{Re } z > \frac{1}{2}$ , the power series representing  $Ch(z)$  converging for  $|z - z^2| < \frac{1}{2}$ . This formula for  $Ch(e)$  as an element of the algebra of differential forms on the algebra of differential forms on the cotangent bundle with modified multiplication) was used already by Fedosov [Fed] in his version of the index theorem. It was also known to A. Connes (private communication) and is more natural and useful than the one given in [Co2].

Let us turn now to the case of an invertible element  $u$  of the unital algebra  $A$ . To simplify the formulas for  $ch_{2n+1}(u)$ ,  $Ch(u)$  we work in the completions of the reduced algebras  $QA' = A * A$ ,  $(\Omega^1 QA)', RA'$ ,  $(\Omega^1 QA)'$  with additional relations  $q(1) = d(1) = \omega(1, x) = \omega(x, 1) = 0$ . Again we easily can define (this time even)

elements  $ch_{2n+1}(u) \in (\varrho A')^n$  as  $\omega(u, u^{-1})^n - \omega(u^{-1}, u)^n$  such that  $d(ch_{2n+1}(u))$  is a sum of supercommutators in  $(\Omega' RA)'$  and such that  $ch_{2n+1}(u)$  differs from  $ch_1(u) = \omega(u, u^{-1}) - \omega(u^{-1}, u)$  by a sum of supercommutators.

The construction of the odd element  $Ch(u)$  is more involved and we have to pass to  $2 \times 2$ -matrices. Put

$$y = \begin{pmatrix} 0 & u^{-1}qu \\ \bar{u}^{-1}qu & 0 \end{pmatrix} \in M_2(\overline{QA}').$$

Using the identity  $u^{-1}q(u)\bar{u}^{-1}q(u) = -q(u^{-1})q(u)$  one finds

$$y^{2n+1} = (-1)^n \begin{pmatrix} 0 & u^{-1}qu(qu^{-1}qu)^n \\ \bar{u}^{-1}qu(qu^{-1}qu)^n & 0 \end{pmatrix}.$$

Given  $c_0 = 1$ , there are unique coefficients  $c_k$  such that  $\partial_{(1)}\varphi = 0$  in the bicomplex dual to 3.B for the chain

$$\varphi = \sum_{n \geq 0} c_{2n} (1 \otimes (u^{-1} \otimes u)^{\otimes n} - 1 \otimes (u \otimes u^{-1})^{\otimes n}) + \sum_{n \geq 0} c_{2n+1} (u^{-1} \otimes u)^{\otimes n+1}$$

and one computes

$$c_{2n} = \frac{2^{2n}(n!)^2}{(2n+1)!}, \quad c_{2n+1} = c_{2n}.$$

Thus, putting

$$f(u) = \sum_{n \geq 0} c_{2n} q(1) ((qu^{-1}qu)^n - (qu qu^{-1})^n) + \sum_{n \geq 0} c_{2n+1} pu^{-1}qu(qu^{-1}qu)^n,$$

$d(f(u))$  is a sum of supercommutators and the first term is 0 if we impose  $q(1) = 0$ . To express  $f(u)$  by a known power series one needs a square root for  $q(u^{-1})q(u)$ , which is given by  $y$ , and we put

$$Ch(u) = \sum_{n \geq 0} (-1)^n c_{2n+1} y^{2n+1}.$$

Given an odd trace  $T$  on  $\overline{QA}'$  we extend  $T$  to  $M_2(\overline{QA}')$  by

$$\tilde{T}\left(\begin{pmatrix} a & b \\ c & d \end{pmatrix}\right) = T(b+c).$$

We then have  $\tilde{T}(Ch(u)) = T(f(u))$  where this number is the product  $\langle [\varphi], [T] \rangle$  of the class of the cycle  $\varphi$  associated with  $u$  and the cyclic cohomology class  $[T]$  associated with  $T$ .

Since the power series  $\sum_{k \geq 0} (-1)^k \frac{2^{2k}(k!)^2}{(2k+1)!} z^{2k+1}$  represents the function  $\frac{\ln(z+\sqrt{1+z^2})}{\sqrt{1+z^2}}$  we formally have

$$Ch(u) = \frac{\ln(y + \sqrt{1+y^2})}{\sqrt{1+y^2}}.$$

To understand the meaning of this formula, suppose that  $\bar{u}$  differs from  $u$  by an exponential:  $\bar{u} = ue^h$ . Then

$$y^2 = \begin{pmatrix} \sinh^2(\frac{h}{2}) & 0 \\ 0 & \sinh^2(\frac{h}{2}) \end{pmatrix}$$

$$\ln(y + \sqrt{1 + y^2}) = \begin{pmatrix} 0 & -\frac{h}{2}e^{\frac{h}{2}} \\ -\frac{h}{2}e^{-\frac{h}{2}} & 0 \end{pmatrix}.$$

This gives

$$\tilde{T}(Ch(u)) = T \left( \frac{-\frac{h}{2}(e^{-\frac{h}{2}} + e^{\frac{h}{2}})}{\cosh(\frac{h}{2})} \right) = -T(h).$$

## References

- [BDF] Brown, L., Douglas, R., Fillmore, P.: Extensions of  $C^*$ -algebras and  $K$ -homology. Ann. Math. (2) **105** (1977) 265–324
- [Co1] Connes, A.: Non commutative differential geometry. Publ. Math. IHES **62** (1985) 257–360
- [Co2] Connes, A.: Entire cyclic cohomology of Banach algebras and characters of  $\Theta$ -summable Fredholm modules.  $K$ -Theory **1** (1988) 519–548
- [Co-Cu] Connes, A., Cuntz, J.: Quasi-homomorphismes, cohomologie cyclique et positivité. Commun. Math. Phys. **114** (1988) 515–526
- [CGM] Connes, A., Gromov, M., Moscovici, H.: Novikov’s conjecture and almost flat vector bundles. C.R. Acad. Sci. Paris, Sér. I **310** (1990) 273–277
- [Cu1] Cuntz, J.: A new look at  $KK$ -theory.  $K$ -theory **1** (1987) 31–51
- [Cu2] Cuntz, J.: Universal extensions and cyclic cohomology. C.R. Acad. Sci. Paris, Sér. I **309** (1989) 5–8
- [Cu-Qu] Cuntz, J., Quillen, D. (In preparation)
- [Fed] Fedosov: Analytic formulas for the index of an elliptic operator. Trudi Mosk. Math. **30** (1974) 159–240 (in Russian)
- [Go] Goodwillie, R.: Cyclic homology, derivations and the free loopspace. Topology **24** (1985) 187–215
- [Lo-Qu] Loday, J.-L., Quillen, D.: Cyclic homology and the Lie algebra homology of matrices. Comment. Math. Helv. **59** (1984) 565–591
- [Ka] Karoubi, M.: Homologie cyclique et  $K$ -théorie. Astérisque **149** (1987)
- [Kas] Kasparov, G.G.: The operator  $K$ -functor and extensions of  $C^*$ -algebras. Izv. Akad. Nauk SSSR, Ser. Mat. **44** (1980) 571–636
- [Qu] Quillen, D.: Algebra cochains and cyclic cohomology. Publ. Math. IHES **68** (1989) 139–174
- [Ts] Tsygan, B.L.: Homology of matrix algebras over rings and Hochschild homology. Uspekhi Mat. Nauk **38** 2 (1983) 217–218
- [Zek] Zekri, R.: A new description of Kasparov’s theory of  $C^*$ -algebra extensions. J. Funct. Anal. **84** (1989) 441–471

# ***K*-Theory for Groups Acting on Trees**

*Michael V. Pimsner*

Mathematisches Institut der Universität Heidelberg, Im Neuenheimer Feld 288  
W-6900 Heidelberg, Fed. Rep. of Germany

*K*-theory is by now such an established subject that it hardly needs any introduction. Recall however that if  $A$  is an associative ring with unit, Grothendieck defined  $K_0(A)$  as the group generated by the semigroup of finitely generated, projective  $A$ -modules. Depending on the nature of the ring  $A$  one gets entirely different theories, one of the most vivid differences appearing for example when trying to define the higher  $K$ -groups.

One of the first achievements of *K*-theory, due to Atiyah, was in topology, in the case  $A = C(X)$ , the algebra of continuous complex valued functions on a compact topological space  $X$ . In this case the finitely generated projective  $A$ -modules are precisely the locally trivial vector bundles, so the whole theory admits a purely topological description. Due to Bott periodicity one gets a  $\mathbb{Z}/2$  graded generalized cohomology theory  $K_*(A)$  (Atiyah 1967). This is a consequence of the topology (and of course of the fact that the ground field are the complex numbers), so it is not really surprising that one gets more or less the same theory for complex Banach algebras (see e.g. Karoubi 1978).

However if  $A = C(X)$  the theory is much richer, since along with *K*-theory comes another dual theory  $K^*(A)$ , called *K*-homology, that has deep connections with elliptic operators (Atiyah 1970), and that may be described either in terms of extensions as (Brown Douglas and Fillmore 1977) did it, or in terms of abstract elliptic operators as (Kasparov 1975) did it. The relevant fact for such a theory to exist is the fact that  $A$  is a  $C^*$ -algebra. Moreover in this case the two dual theories can be glued together. This was first done in (Pimsner, Popa, and Voiculescu 1979 and 1980) in the particular case when one variable is still commutative (and finite dimensional), motivated by the study of the homotopy invariance of extensions, and soon later (Kasparov 1980) came with what looks now to be the final theory. This is a functor of two (not necessarily commutative)  $C^*$ -algebras  $A$  and  $B$ , (which we will assume from now on for simplicity to be separable), with values in the  $\mathbb{Z}/2$  graded abelian groups, denoted  $KK_*(A, B)$ , such that

$$K_*(A) = KK_*(\mathbb{C}, A)$$

and

$$K^*(A) = KK_*(A, \mathbb{C})$$

( $\mathbb{C}$  stands as usual for the complex numbers). The most important feature of this theory is the product introduced by Kasparov, which is the analogue of the various cup, cap, slant... products in cohomology theory. In its simplest form it

may be described as a map

$$KK_*(A, B) \times KK_*(B, C) \longrightarrow KK_*(A, C)$$

denoted  $(x, y) \rightarrow x \otimes_B y$ .

The central theorem of  $KK$ -theory gives the construction of the product and asserts its associativity. For example the pairing between  $K$ -theory and  $K$ -homology is just the Kasparov product

$$KK(\mathbb{C}, A) \otimes KK(A, \mathbb{C}) \longrightarrow KK(\mathbb{C}, \mathbb{C})$$

since the latter is isomorphic to the additive group of integers  $\mathbb{Z}$ .

Another important consequence of the Kasparov product is the fact that one can regard the elements of  $KK_*(A, B)$  as natural maps from  $K_*(A)$  to  $K_*(B)$ ,

$$KK_*(\mathbb{C}, A) \times KK_*(A, B) \longrightarrow KK_*(\mathbb{C}, B)$$

with composition given by the product. So even if one is interested only in  $K$ -theory, it is unnatural to discard  $KK$ -theory, since most of the interesting maps are not given by \*-homomorphisms from  $A$  to  $B$  but by elements from  $KK_*(A, B)$ . An improvement of this aspect of  $KK$ -theory was recently proposed by Connes and Higson (1990).

Another consequence of the Kasparov product is the fact that  $KK_*(A, A)$  becomes a (graded) ring with unit  $1_A$ , given by the class of the identity map from  $A$  to  $A$ , and that  $KK_*(A, B)$  has a left  $KK_*(A, A)$  (and a right  $KK_*(B, B)$ ) module structure.

It should be stressed at this point that to show that the product of two given elements equals a third one is often a deep theorem. It is enough to mention for example that the Atiyah-Singer theorem (Atiyah and Singer 1968) can be put into this form (Kasparov 1975, Connes and Skandalis 1984, Baum and Douglas 1982).

Of course  $KK$ -theory did not appear just for the sake of generalizing topological  $K$ -theory from spaces to general  $C^*$ -algebras. Already in the seventies  $K_*$  and  $K^*$  proved to be powerful invariants for certain classes of  $C^*$ -algebras (Elliott 1976, Pimsner and Popa 1978). One of the turning points in applying  $K$ -theory to  $C^*$ -algebras was probably the discovery (Pimsner and Voiculescu 1980), that the irrational rotation  $C^*$ -algebras  $A_\theta$  are essentially nonisomorphic for different  $\theta$ 's. These are the  $C^*$ -algebras generated in  $L^2$  of the unit circle by multiplication operators with continuous functions and by the rotation operator by the angle  $2\pi\theta$ . Their importance lies in the fact that they are the most simple nontrivial  $C^*$ -algebras of several classes of  $C^*$ -algebras. One can look at them as crossed products by  $\mathbb{Z}$ , as the group  $C^*$ -algebra of  $\mathbb{Z}^2$  with a 2-cocycle (see e.g. Pedersen 1979) or as stably isomorphic to the  $C^*$ -algebra of the Kronecker foliation (Connes 1982). The computation of their  $K_*$ - and  $K^*$ -theories was done in (Pimsner and Voiculescu 1980) by proving a general six terms exact sequence for every crossed product by  $\mathbb{Z}$ . In the particular case of the irrational rotation algebras one gets:

**Theorem** (Pimsner and Voiculescu 1980).

$$K_0(A_\theta) \simeq \mathbb{Z}^2$$

with generators [1] and [R], the class of the identity and respectively that of the Rieffel projection.

$$K_1(A_\theta) \simeq \mathbb{Z}^2$$

with generators [u] and [v], the unitaries corresponding to the rotation and respectively to multiplication by z.

Surprisingly enough the same techniques could be used to solve the conjecture of Kaplansky-Kadison, that the reduced  $C^*$ -algebra of the free group on two generators (which was proved to be simple by (Powers 1975) ) had no nontrivial (i.e. different from 0 and 1) projection. This was done (Pimsner Voiculescu 1982), again by proving a six terms exact sequence for reduced crossed products by free groups on any number of generators (not just one). In particular one gets:

**Theorem** (Pimsner and Voiculescu 1982).

$$K_0(C_r^*(\mathbb{F}_n)) \simeq \mathbb{Z},$$

with generator [1], the class of 1.

$$K_1(C_r^*(\mathbb{F}_n)) \simeq \mathbb{Z}^n$$

with generators  $[g_1], \dots, [g_n]$ , the classes of the generators of  $\mathbb{F}_n$ .

As a corollary of this theorem one gets also that the reduced  $C^*$ -algebra of free groups are nonisomorphic if the number of generators are distinct. This may be just one hint of why the corresponding question for the von Neumann algebras of the free groups is so difficult.

Another motivation for studying the K-theory of crossed products by discrete groups came from topology. I won't pursue this any further, I will just mention that some knowledge of the K-theory of  $\pi_1(M)$  (the fundamental group of the manifold M) gives a positive answer to the Novikov conjecture about the homotopy invariance of higher signatures for manifolds with prescribed homotopy group (Kasparov 1988).

It is quite clear that the K-theory of crossed products is important both for the study of  $C^*$ -algebras and for geometric applications. The basic tool for this problem is the equivariant  $KK^G$ -theory developed by Kasparov. Unlike the compact group case that can be obtained quite easily by taking only equivariant elements in the definition of the KK-groups, the general noncompact case is much more difficult and shows one of the advantages of KK-theory. Relevant for our purposes is the induction map

$$j : KK_*^G(A, B) \longrightarrow KK_*(A \rtimes G, B \rtimes G)$$

that commutes with the Kasparov product and sends the unit of  $KK_*^G(A, A)$  to the unit of  $KK_*(A \rtimes G, A \rtimes G)$ . In particular it transforms  $KK^G$ -equivalence into KK-equivalence, where we have the following natural definition.

**Definition.** The  $C^*$ -algebras A and B are KK-equivalent if there exist  $\alpha \in KK(A, B)$  and  $\beta \in KK(B, A)$  such that  $\alpha \otimes_B \beta = 1_A$  and  $\beta \otimes_A \alpha = 1_B$ .

(The  $KK^G$ -equivalence is defined in the same way using  $KK^G$ -theory). Up to now there are two ways of studying the K-theory of crossed products by discrete

groups. One way, emphasized by Connes and Kasparov, takes the point of view that it is easier to deal with Lie groups, which then in turn will provide results for all their discrete subgroups. The second one treats the discrete group on its own right. Both methods are of geometric nature the main ingredient being the study of “spaces” on which the groups act naturally.

To show how this works let me describe, at least at a formal level, the Connes-Kasparov conjecture. Let  $G$  be a connected Lie group,  $K$  its maximal compact subgroup and  $X = G/K$  the corresponding homogeneous space of dimension  $n$ . The class of the Dirac operator defines an element  $\alpha \in KK_n^G(C_0(X), \mathbb{C})$ . The needed  $\beta \in KK_n^G(\mathbb{C}, C_0(X))$  is the element defined by Kasparov (1973) and Mishchenko (1973) out of the negative curvature of the homogeneous space  $X$ .

**Theorem (Kasparov).**

1.  $\alpha \otimes_{\mathbb{C}} \beta = 1_{C_0(X)} \in KK_0^G(C_0(X), C_0(X))$
2.  $\beta \otimes_{C_0(X)} \alpha = \gamma^G \in KK_0^G(\mathbb{C}, \mathbb{C})$

is an idempotent.

So if one can prove that  $\gamma^G = 1_{\mathbb{C}}$  for the group  $G$ , then  $C_0(X)$  is  $KK_n^G$ -equivalent to  $\mathbb{C}$  and taking crossed products one gets that  $C^*(G)$  is  $KK_n$ -equivalent to  $C_0(X) \rtimes G$ . Since this latter  $C^*$ -algebra is strong Morita equivalent (in the sense of Rieffel (1974); this is a particular, more easy,  $KK$ -equivalence) to  $C^*(K)$  one gets in particular

$$K_*(C^*(G)) \simeq K_{*+n}(C^*(K)),$$

the isomorphism being natural.

Moreover if  $\Gamma$  is a discrete subgroup of the connected Lie group  $G$ , with  $\gamma^G = 1_{\mathbb{C}}$ , then the same holds:  $C_0(X)$  being  $KK_n^\Gamma$ -equivalent to  $\mathbb{C}$ , one gets by taking crossed products by  $\Gamma$  that  $C^*(\Gamma) \simeq_{KK_n} C_0(X) \rtimes \Gamma$  and again the latter  $C^*$ -algebra is strongly Morita equivalent to  $C_0(X/\Gamma) \rtimes K$ . Still better, we can start with  $C_0(X) \otimes A \simeq_{KK_n^\Gamma} A$  to get the  $K$ -theory of  $A \rtimes \Gamma$ , for any  $C^*$ -algebra on which  $\Gamma$  acts. However this depends on  $\gamma^G = 1_{\mathbb{C}}$  (or at least  $\gamma^\Gamma = 1_{\mathbb{C}}$  for the discrete subgroup  $\Gamma$ ). This is known for a large class of groups, including the amenable groups,  $SO(n, 1)$  (Kasparov 1984),  $SU(n, 1)$  (Julg and Kasparov 1990). It is also known to fail for groups having property  $T$  of Kazhdan. Note however that it would be enough to have  $C^*(\Gamma) \simeq_{KK_n} C_0(X) \rtimes \Gamma$ . In order to study this question Cuntz (1983) introduced the notion of  $K$ -amenability. But even this is not true in general as has been shown in a remarkable paper by Skandalis (1988).

Thus in general we can compute only what is called the “ $\gamma$ -part” of the  $K$ -groups, which is a direct summand of what we really want, since  $\gamma^2 = \gamma$ . If  $\gamma^\Gamma = 1_{\mathbb{C}}$  (for the discrete subgroup of the Lie group  $G$ ), we shall say that we know the  $K$ -theory for the group  $\Gamma$ . Note that this is really much more than knowing the  $K$ -groups of the crossed products.

Let us stop here this very incomplete presentation of the Lie group case and turn now to the case of groups that act on some oriented tree  $X$ , that is on a one-dimensional simply connected simplicial complex  $X = (X^1, X^0)$ , where  $X^1$  will be the set of edges and  $X^0$  the set of vertices (see Serre 1977). Since the tree is oriented there are two maps  $t, o : X^1 \rightarrow X^0$ , which are thought to be the

terminus and respectively the origin of the edge. An action of the group  $G$  on  $X$  is then an action of  $G$  on both  $X^1$  and  $X^0$ , such that the maps  $t$  and  $o$  are  $G$ -equivariant.

The case of the discrete groups that act on some tree is well understood due to the Bass-Serre theory of graphs of groups (see Serre 1977). These consist of an oriented graph  $\Sigma = (\Sigma^1, \Sigma^0)$  and of a collection  $\{G_y\}_{y \in \Sigma^1}$  and  $\{G_P\}_{P \in \Sigma^0}$  of groups, together with embeddings  $t : G_y \rightarrow G_{t(y)}$  and  $o : G_y \rightarrow G_{o(y)}$  for every  $y \in \Sigma^1$ . Out of a graph of groups one gets a group by constructing the so called “fundamental group of the graph of groups”, that acts on the “universal cover of the graph of groups” which is a tree. The inverse construction is more straightforward:  $\Sigma$  is simply the orbit (homogeneous) space  $G \backslash X$ , while the  $G_y$ 's (respectively the  $G_P$ 's) are the stabilizer subgroups of the edges (resp. of the vertices) modulo some automorphism coming from the choice of a representative. The simplest graphs of groups, having only one edge, yield well known constructions in the theory of discrete groups, namely amalgamated products (if there are 2 vertices) and *HNN*-extensions (if there is only one vertex).

The nondiscrete groups that act on some tree are as interesting as the discrete ones (Serre 1977). The most interesting seem to be the reductive groups over local fields with one dimensional Bruhat-Tits building (see Tits 1979), e.g.  $SL_2(\mathbb{Q}_p)$ .

Since the computation of the  $K$ -groups of the reduced crossed products by free groups (Pimsner and Voiculescu 1982), partial results on the  $K$ -groups of free and amalgamated products and of *HNN*-extensions of groups have been obtained. Thus Lance (1983) introduced condition  $A$ , in order to use the methods of (Pimsner Voiculescu 1982) to compute the  $K$ -groups of the reduced  $C^*$ -algebra of certain free products of groups. This has been extended by Natsume (1985) to certain amalgamated products and finally Anderson and Paschke (1986) combined the above results with those of (Pimsner Voiculescu 1980) to get results for the  $K$ -groups of the reduced  $C^*$ -algebra of certain *HNN*-extensions. However the action of the group  $G$  on the tree  $X$  allows us to get much more, namely the knowledge of the *K-theory for the group G*. I will describe the original *Toeplitz Extension* approach of (Pimsner 1986), which is the generalisation of the methods of (Pimsner and Voiculescu 1980, 1982) to the general tree case. It is based on a rough analysis of actions of groups on trees due to Julg and Valette (1984). To this end one fixes a vertex  $O$ , called origin, and one denotes by  $\chi_O$  the set of all edges that point to  $O$ . By adding to the continuous functions on the discrete space  $X^1$  that vanish at infinity the characteristic function of  $\chi_O$ , one gets one point at infinity that is fixed by the action of  $G$ . Denoting by  $X_+^1$  the space thus constructed one gets the following  $G$ -equivariant exact sequence:

$$0 \longrightarrow C_0(X^1) \longrightarrow C_0(X_+^1) \longrightarrow \mathbb{C} \longrightarrow 0$$

which we call the ( $G$ -equivariant) Toeplitz extension. The name comes from the case of  $\mathbb{Z}$  acting on its obvious tree, since after taking crossed products one gets the exact sequence

$$0 \longrightarrow K(L^2(\mathbb{T})) \longrightarrow T \longrightarrow C(\mathbb{T}) \longrightarrow 0$$

and where  $\chi_O$  corresponds to the Hardy projection onto  $H^2(\mathbb{T})$ . With these notations in mind we have the following theorem.

**Theorem.** *The  $C^*$ -algebra  $C_0(X_+^1)$  is  $KK_0^G$ -equivalent to the  $C^*$ -algebra  $C_0(X^0)$ .*

This really means that there are elements  $\alpha \in KK_0^G(C_0(X^0), C_0(X_+^1))$  and  $\beta \in KK_0^G(C_0(X_+^1), C_0(X^0))$  and that one is the inverse of the other in  $KK^G$ . This theorem together with the Toeplitz extension gives the  $K$ -theory for groups acting on trees. For example to compute the  $K$ -groups of crossed products, one tensors the Toeplitz extension with the  $C^*$ -algebra  $A$ , takes crossed products with  $G$  and applies the  $K$ -theory functor. Replacing in the periodic six-terms exact sequence the  $KK$ -equivalent terms one gets the following theorem

**Theorem** (Pimsner 1986). *Let  $G$  be a second countable group that acts on some oriented tree  $X$  and on the  $C^*$ -algebra  $A$ . The following six terms periodic sequence is exact:*

$$\begin{array}{ccccccc} \oplus_{y \in \Sigma^1} K_0(A \rtimes G_y) & \xrightarrow{\sigma} & \oplus_{P \in \Sigma^0} K_0(A \rtimes G_P) & \xrightarrow{i} & K_0(A \rtimes G) \\ \uparrow \delta & & & & \downarrow \delta \\ K_1(A \rtimes G) & \xleftarrow{i} & \oplus_{P \in \Sigma^0} K_1(A \rtimes G_P) & \xleftarrow{\sigma} & \oplus_{y \in \Sigma^1} K_1(A \rtimes G_y) \end{array}$$

where  $\sigma = t_* - o_*$  is the obvious map given by the difference of the terminus and origin maps,  $i$  is the map given by the inclusion maps  $G_P \rightarrow G$  and  $\delta$  are the maps induced by the boundary maps associated to the Toeplitz extension.

This is the generalisation of the six terms exact sequence for the free groups (Pimsner and Voiculescu 1982). Note that the above result is expressed in terms of the graph of groups associated to the action of  $G$  on  $X$ , so that one gets explicit exact sequences that express the  $K$ -groups of the group by the  $K$ -groups of the corresponding graph of groups. In particular one gets exact sequences for amalgamated products and for HNN-extensions of groups. Of course the same reasoning gives also the  $K^*$ -groups and more general the  $KK^*$ -groups of crossed products by groups acting on trees.

There is one technical point concerning crossed products which we did not mention at all, namely the difference between the reduced and the full ones. One reason is that the exact sequence exists for each of them, provided we fix one of the cross-norms. This has however an interesting corollary, which is a generalization of a result of Julg and Valette (1984)

**Corollary.** *If the countable discrete group  $G$  acts on some tree  $X$ , then  $G$  is  $K$ -amenable if and only if every stabilizer is  $K$ -amenable*

As we already mentioned, the Toeplitz extension together with the  $KK^G$ -equivalence of the Toeplitz algebra with  $C_0(X^0)$ , give us in fact more than just the above exact sequences. To illustrate this let us consider  $\Gamma = \Gamma_1 \times \dots \times \Gamma_n$  the direct product of groups, each of them acting on some tree  $X_i$ . Then one still can express the  $KK$ -groups of the crossed product by  $\Gamma$  in terms of the  $KK$ -groups of the crossed products by the stabilizer subgroups. This time one gets a *spectral sequence* associated to the simplicial decomposition of the  $n$ -dimensional simplicial complex  $X = X_1 \times \dots \times X_n$ . This is due to the following obvious consequence of the theorem

$$\oplus_{i_1, \dots, i_n} C_0(\dots \times X_{i_1,+}^1 \times \dots \times X_{i_p,+}^1 \times \dots) \simeq_{KK^{\Gamma}} C_0(X^{n-p})$$

where the sum in the left hand side of the formula is taken over all direct products that have on exactly  $p$  positions the edges with the point at infinity added and just the edges on the other ones.

This remark brings us to the last topic, namely the generalization to groups acting on buildings due to Kasparov and Skandalis (1989). Their approach makes the analogy with the Lie groups and with the Connes Kasparov conjecture even more transparent. Regarding the building  $X$  as a non Hausdorff manifold of negative sectional curvature they construct the  $C^*$ -algebra  $\mathcal{A}_X$ , which is the right analogue of the  $C^*$ -algebra of continuous functions on  $X$  and under the additional assumption that the building is locally finite, the “Dirac” element  $\alpha \in KK_n^\Gamma(\mathcal{A}_X, \mathbb{C})$  and the “dual Dirac” element  $\beta \in KK_n^\Gamma(\mathbb{C}, \mathcal{A}_X)$ . Moreover they prove the following theorem.

**Theorem (Kasparov Skandalis).**

1.  $\alpha \otimes_{\mathbb{C}} \beta = 1_{\mathcal{A}_X}$
2.  $\beta \otimes_{\mathcal{A}_X} \alpha = \gamma^\Gamma$

is an idempotent.

The idempotent  $\gamma^\Gamma$  has been previously constructed by Julg and Valette (1988). They also showed (1984) that it equals  $1_{\mathbb{C}}$  in the tree case. As in the case of Lie groups, if  $\gamma^\Gamma = 1_{\mathbb{C}}$ , we know the  $K$ -theory for  $\Gamma$ . In this case the crossed products by  $\Gamma$  are expressed in terms of a spectral sequence of the crossed products by the stabilizer subgroups. In general only the  $\gamma$ -part is known. This brings the study of the  $K$ -theory of groups acting on buildings at the same level as that of connected Lie groups. On one hand there is the  $\gamma$ -obstruction, on the other hand we know the positive answer to Novikov’s conjecture for such groups (Kasparov 1988, Kasparov and Skandalis 1989).

However the one dimensional case, that is the case of a tree, is special. For first of all one does not need the locally finiteness of the tree, and second in this case  $\gamma = 1$  for geometric reasons.

## References

- Anderson, J., Paschke, W. (1986): The  $K$ -theory of the reduced  $C^*$ -algebra of an HNN-extension. *J. Operator Theory* **16** (1986) 165–187
- Atiyah, M.F. (1967):  $K$ -theory. Benjamin, New York Amsterdam
- Atiyah, M.F. (1970): Global theory of elliptic operators. In: Proc. Intern. Conf. on Functional Analysis and Related Topics, Univ. of Tokyo Press, Tokyo, pp. 21–30
- Atiyah, M.F., Singer, I.M. (1968): The index of elliptic operators III. *Ann. Math.* **87**, 546–604
- Baum, P., Douglas, R.G. (1982):  $K$ -homology and index theory. In: Operator Algebras and Applications, Proc. Symp. Pure Math. Amer. Math. Soc. **38**, part 1, 117–173
- Brown, L.G., Douglas, R.G., Fillmore, P.A. (1977) : Extensions of  $C^*$ -algebras and  $K$ -homology. *Ann. Math.* **105**, no.2, 265–324
- Connes, A. (1982): A Survey of foliations and operator algebras. In: Operator Algebras and Applications, Proc. Symp. Pure Math. Amer. Math. Soc. **38**, part 1, 521–628
- Connes, A., Higson, N. (1990): Déformations, morphismes asymptotiques et  $K$ -théorie bivariante. Preprint

- Connes, A., Skandalis, G. (1984): The longitudinal index theorem for foliations. *Publ. Res. Inst. Math. Sci. Kyoto Univ.* **20**, 1139–1183
- Cuntz, J. (1983): K-theoretic amenability for discrete groups. *J. Reine Angew. Math.* **344**, 180–195
- Elliott, G.A. (1976): On the classification of inductive limits of sequences of semisimple finite dimensional algebras. *J. Algebra* **38**, 29–44
- Julg, P., Valette, A. (1984): K-amenable for  $SL_2(\mathbb{Q}_p)$  and the action on the associated tree. *J. Funct. Anal.* **58**, 194–215
- Julg, P., Valette, A. (1988): Fredholm modules associated to Bruhat-Tits buildings. Preprint
- Julg, P., Kasparov, G.G. (1990): L’anneau  $KK_G(\mathbb{C}, \mathbb{C})$  pour  $G = SU(n, 1)$ . Preprint
- Karoubi, M. (1978): K-theory. An Introduction. Springer, Berlin Heidelberg New York
- Kasparov, G.G. (1973): The generalized index of elliptic operators. *Funkt. Anal. Prilozhen.* **7**, no. 3, 82–83 (in Russian)
- Kasparov, G.G. (1975): Topological invariants of elliptic operators. In: K-homology. *Izv. Akad. Nauk SSSR, Ser. Mat.* **39**, 796–838
- Kasparov, G.G. (1980): The operator K-functor and extensions of  $C^*$ -algebras. *Izv. Akad. Nauk SSSR, Ser. Mat.* **44**, 571–636
- Kasparov, G.G. (1984): Lorentz groups: K-theory of unitary representations and crossed products. *Dokl. Akad. Nauk USSR*, **275**, 541–545 (in Russian)
- Kasparov, G.G. (1988): Equivariant KK-theory and the Novikov conjecture. *Invent. math.* **91**, 147–201
- Kasparov, G.G., Skandalis, G. (1989): Groups acting on buildings, operator K-theory, and Novikov’s conjecture. Preprint
- Lance, E.C. (1983): K-theory for certain group  $C^*$ -algebras. *Acta. Math.* **151**, 209–230
- Mishchenko, A.S. (1973): Infinite dimensional representations of discrete groups and homotopic invariants of nonsimply connected manifolds. *Usp. Mat. Nauk* **28**, no. 2, 239–240 (in Russian)
- Natsume, T. (1985): On  $K_*(C^*(SL_2(\mathbb{Z})))$ . *J. Operator Theory* **13**, 103–118
- Pedersen, G.K. (1979):  $C^*$ -algebras and their automorphism groups. Academic Press, New York London
- Pimsner, M.V. (1986): KK-groups of crossed products by groups acting on trees. *Invent. math.* **86**, 603–634
- Pimsner, M.V. (1987): Cocycles on trees. *J. Operator Theory* **17**, 121–128
- Pimsner, M.V., Popa, S.T. (1978): The Ext-groups of some  $C^*$ -algebras considered by J. Cuntz. *Rev. Roum. de Math. Pures et Appl.*, Tome XXII 7, 1069–1076
- Pimsner, M.V., Popa, S.T., Voiculescu, D.V. (1979): Homogeneous  $C^*$ -extensions of  $C(X) \otimes K(H)$ , part 1. *J. Operator Theory* **1**, no. 1, 55–108
- Pimsner, M.V., Popa, S.T., Voiculescu, D.V. (1980): Homogeneous  $C^*$ -extensions of  $C(X) \otimes K(H)$ , part 2. *J. Operator Theory* **4**, no. 2, 211–249
- Pimsner, M.V., Voiculescu, D.V. (1980): Exact sequences for K-groups and Ext-groups of certain cross-product  $C^*$ -algebras. *J. Operator Theory* **4**, no. 1, 93–118
- Pimsner, M.V., Voiculescu, D.V. (1982): K-groups of reduced crossed products by free groups. *J. Operator Theory* **8**, no. 1, 131–156
- Powers, R.T. (1975): Simplicity of the  $C^*$ -algebra associated with the free group on two generators. *Duke Math. J.* **42**, 151–156
- Rieffel, M.A. (1974): Induced representations of  $C^*$ -algebras. *Adv. Math.* **13**, 176–257
- Serre, J.P. (1977): Arbres, amalgames,  $SL_2$ . Astérisque **46**, (Soc. Math. France)
- Skandalis, G. (1988): Une Notion de nucléarité en K-théorie. *K-Theory* **1**, 549–573
- Tits, J. (1979): Reductive groups over local fields. *Proc. Symp. Pure Math.* **33**, 29–69

# Subfactors and Classification in von Neumann Algebras

*Sorin Teodor Popa*

University of California, Los Angeles, CA 90024, USA

## 0. Introduction

While initially introduced to study quantum mechanics and representations of groups, in recent years von Neumann algebras started to play a major role in many areas of mathematics. The class of von Neumann algebras that proved to be more important and of most physical significance are the ones that can be approximated by finite dimensional algebras, called hyperfinite. By 1985, as a result of 45 years of work initiated by Murray and von Neumann and culminating with the work of Connes and Connes-Haagerup, hyperfinite algebras were completely classified. The fundamental steps in accomplishing this proved to be the uniqueness of the hyperfinite type  $\text{II}_\infty$  factor (the hyperfinite factor with an infinite trace) and the classification of its automorphisms, both settled by Connes. Starting with his work on automorphisms, the problem of classifying actions of more general groups on the hyperfinite factors became an important trend of research in von Neumann algebras ([J1, Oc1, JT]).

In 1983, motivated by his study of actions of finite groups on hyperfinite von Neumann algebras, Jones initiated a Galois theory for von Neumann algebras by studying pairs of factors  $N \subset M$  with finite index  $[M : N] < \infty$ , i.e. with  $M$  a finite  $N$  module. He proved the striking result that if the index of  $N \subset M$  is less than 4, then it has to be equal to the square norm of a matrix with nonnegative integer entries, and must be of the form  $4 \cos^2 \frac{\pi}{n}$  for some  $n \geq 3$ . Jones' results and ideas gave a new insight into the theory of operator algebras. It also brought together many other subjects and had an unexpected and far reaching impact in a number of fields of research such as statistical mechanics, quantum field theory, knot theory. Methods and results in either of these subjects proved to be inspiring for the others. Most of the intrinsic problems from the theory of subfactors, such as their classification, the construction of examples, the characterization of the values the index may take, seem to have physical significance. We will describe in this article the solution to some of these problems. We will present a classification result for a certain class of subfactors, called strongly amenable subfactors, a class that contains all subfactors of index  $\leq 4$ . We will show that the index of an irreducible subfactor  $N \subset M$  of the hyperfinite factor must always be the square norm of a (possibly infinite) matrix of nonnegative integers. In particular this

shows the existence of a gap in the set of indices, from 4 to 4.026.... Finally we will show that in striking contrast with the hyperfinite case, in the nonhyperfinite case any value  $> 4$  may appear as the index of an irreducible subfactor. We will also explain how the theory of subfactors provides the natural framework and the necessary techniques for a unified approach to both subfactors problems and the classical operator algebra problems mentioned before: classification of hyperfinite von Neumann algebras and of their automorphisms.

## 1. Index for Subfactors

Let  $M$  be a type  $\text{II}_1$  factor with a normalized trace  $\tau$  and  $N \subset M$  a subfactor. The Jones' index of  $N$  in  $M$ , denoted by  $[M : N]$ , is defined as  $\dim_N L^2(M)$ , the Murray-von Neumann coupling constant of  $N$  in its representation on  $L^2(M)$ ,  $L^2(M)$  being the completion of  $M$  in the Hilbert norm  $\|x\|_2 = \tau(x^*x)^{1/2}$ ,  $x \in M$ . The definition of  $[M : N]$  is in fact independent of the Hilbert space on which  $N \subset M$  act, as  $[M : N] = \dim_N \mathcal{H}/\dim_M \mathcal{H}$  for any  $M$ -Hilbert module  $\mathcal{H}$ . But, surprisingly enough, although each of  $\dim_N \mathcal{H}$ ,  $\dim_M \mathcal{H}$  may take any value from 0 to  $\infty$  (a remarkable fact in itself!), Jones proved in [J2] that  $[M : N]$  may only take the values  $\{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\} \cup [4, \infty)$ . The key idea in his proof is the so called basic construction. It later proved to be the fundamental construction of the theory. This construction associates to  $N \subset M$  a new pair of factors  $M \subset M_1$ , with the same index  $[M_1 : M] = [M : N]$ , and with a projection  $e_1 \in M_1$  that generates  $M_1$  together with  $M$  and implements by compression the trace preserving conditional expectation of  $M$  onto  $N$ . By iteration one can thus produce a whole “tower” of factors  $N \subset M \subset M_1 \subset \dots$  together with a sequence of projections  $e_i \in M_i$ ,  $i \geq 1$ , satisfying the remarkable axioms:

- 1.1.1  $[e_i, e_j] = 0$ ,  $|i - j| \geq 2$ .
- 1.1.2  $e_i e_{i \pm 1} e_i = [M : N]^{-1} e_i$ ,  $i \geq 1$ .
- 1.1.3  $\tau(w e_{n+1}) = [M : N]^{-1} \tau(w)$ ,  $w \in \text{Alg}\{1, e_1, \dots, e_n\}$ .

Conditions 1.1.1–1.1.3 is what forces  $s = [M : N]$  to be in the set  $\{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\}$ , if less than 4. Moreover, for each  $s \in \{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\} \cup [4, \infty)$  Jones proved the existence of such a sequence of projections with a trace satisfying 1.1.1–1.1.3. Then  $R = \overline{\text{Alg}\{e_i\}}_{i \geq 1}^w$  is isomorphic to the hyperfinite type  $\text{II}_1$  factor and  $R^s = \overline{\text{Alg}\{e_i\}}_{i \geq 2}^w$  is a subfactor of  $R$  of index  $[R : R^s] = s$ . So one has:

- 1.2 Theorem** [J2]. a)  $[M : N] \in \{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\} \cup [4, \infty)$ .  
b) Given any  $s \in \{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\} \cup [4, \infty)$  there exists a pair of hyperfinite type  $\text{II}_1$  factors  $R^s \subset R$  with index  $s$ .

Beside the above Jones' subfactors, we mention few other important classes of examples:

- 1.3.1 If  $M = M_{n \times n}(N)$  is the algebra of  $n$  by  $n$  matrices over  $N$ , then  $[M : N] = n^2$ .

1.3.2 If  $\sigma$  is a properly outer action of a finite group  $G$  on a type  $\text{II}_1$  factor  $N$  then  $N \subset M = N \rtimes_{\sigma} G$  satisfies  $[M : N] = |G|$ . If  $\mu$  is another action of  $G$  on  $N$  then  $N \subset N \rtimes_{\sigma} G$  is isomorphic to  $N \rtimes_{\mu} G$  if and only if  $\sigma$  is cocycle conjugate to  $\mu$ . The classification of the actions  $\sigma$  is thus equivalent to the classification of the corresponding pairs of algebras, i.e. subfactors.

1.3.3 Let  $G$  be a finitely generated discrete group and fix  $g_1, \dots, g_n$ , some generators. Let  $\sigma$  be an outer action of  $G$  on a type  $\text{II}_1$  factor  $P$  or, more generally, an injective morphism of  $G$  into  $\text{Aut } P / \text{Int } P$ . Let  $M = M_{(n+1) \times (n+1)}(P)$  and  $N_{G,\sigma} = \{x \oplus (\oplus_i \sigma(g_i)(x)) \mid x \in P\} \subset M$  be the “diagonal” subfactor (in general  $\sigma(g_i)$  are lifting automorphisms). Then  $[M : N_{G,\sigma}] = (n+1)^2$ . The isomorphism class of  $N_{G,\sigma} \subset M$  doesn’t depend on unitary perturbations of  $\sigma$ . In fact,  $N_{G,\sigma} \subset M$  is isomorphic to  $N_{G,\mu} \subset M$  iff  $\sigma$  is cocycle conjugate to  $\mu$ , i.e.  $\sigma$  and  $\mu$  are conjugate in  $\text{Aut } P / \text{Int } P$  (for all this see [Po4]).

1.3.4 The Jones subfactors  $R^s \subset R$  in 1.2 come from certain positive Markov traces on the infinite Hecke algebras  $H_{\infty}(q)$  and are related to the Jones’ polynomial invariant for knots. The general form of such a type of subfactor, obtained by investigating all possible positive Markov traces on  $H_{\infty}(q)$  (found by Ocneanu, see [J3]), were constructed by Wenzl [We1], who computed their indices and proved that they are irreducible when  $q$  and the parameters describing the trace correspond to roots of unity.

1.3.5 Let  $G$  be a compact group acting minimally on a type  $\text{II}_1$  factor  $M$  (i.e. such that the fixed point algebra  $M^G$  is irreducible in  $M$ ). Let  $\pi$  be a unitary representation of  $G$  on a finite dimensional space  $V$ . Then  $M^G \subset (M \otimes \text{End } V)^G$  is an inclusion of type  $\text{II}_1$  factors of index  $(\dim V)^2$  which is irreducible if  $\pi$  is. Its basic construction is obtained by tensoring recursively with  $\text{End } V$  and by taking fixed point algebras. This example of subfactors is due to Wassermann [Wa].

1.3.6 Let  $\{A_n \subset B_n\}_{n \geq 1}$  be an increasing sequence of inclusions of finite dimensional algebras, i.e.  $A_n \subset A_{n+1}$ ,  $B_n \subset B_{n+1}$ , with a common trace  $\tau$ . Assume the consecutive steps of these inclusions satisfy the so-called *commuting square condition*:  $E_{A_{n+1}} E_{B_n} = E_{A_n}$ ,  $n \geq 1$ . Then, under suitable conditions on the trace, one has that  $N = \overline{\cup A_n} \subset \overline{\cup B_n} = M$  are factors and  $[M : N] = \lim [B_n : A_n]$ , with  $[B_n : A_n]$  appropriately defined ([PiPo1, PiPo2]). Moreover if  $T_n$  denotes the inclusion matrix for  $A_n \subset B_n$ , and  $\tau$  is extremal in some sense then  $[M : N] = \lim \|T_n\|^2$  ([PiPo2]).

## 2. Classification of Amenable Subfactors

The problem of classifying the subfactors of the hyperfinite type  $\text{II}_1$  factor is of most importance for the theory of subfactors. Moreover, the example 1.3.3 shows that a classical problem such as the classification of actions of groups on factors can be translated into a classification problem for subfactors.

The classification of  $N \subset M$  we are interested in is up to conjugacy of  $N$  by automorphisms of  $M$ . We will present now a result that classifies an important class of subfactors in terms of certain canonically associated combinatorial data.

**2.1 Definition of the Invariant.** Let  $\{M'_1 \cap M_k \subset M' \cap M_k\}_k$  be the sequence of inclusion of finite dimensional algebras of the higher relative commutants in the Jones' tower of factors. Consecutive steps of these inclusions satisfy the commuting square condition 1.3.6. The isomorphism class of this sequence depends only on the isomorphism class of  $N \subset M$ , so it is an invariant for  $N \subset M$ . We call it the *sequence of canonical (or standard) commuting squares* associated to  $N \subset M$ . The finite dimensional inclusions  $M' \cap M_k \subset M' \cap M_{k+1}$ ,  $k \geq 1$ , are all described by a unique irreducible matrix of nonnegative integers  $A$ , called the *standard matrix* of  $N \subset M$ , and respectively by its transpose (for  $k$  even resp. odd), see [GHJ, Oc2, We1, Po2].

The standard matrix  $A$  can alternatively be regarded as a diagram or a graph. As first noted by Jones, for small values of the index ( $< 4$ ) such a graph is necessarily a Coxeter graph of type  $A_n$ ,  $D_n$ ,  $E_6$ ,  $E_7$ ,  $E_8$ . Although the matrices (or graphs)  $A$  capture a great deal of information on the canonical sequence, in general it doesn't uniquely determine it. For the class of subfactors for which  $A$  is finite, a condition introduced by Ocneanu [Oc2] and that he calls *finite depth*, the canonical sequence of commuting squares is uniquely determined by just one of the commuting squares, involving two consecutive steps of the sequence

$$\begin{array}{ccc} M'_1 \cap M_{k+1} & \subset & M' \cap M_{k+1} \\ \cup & & \cup \\ M'_1 \cap M_k & \subset & M' \cap M_k \end{array}$$

with  $k$  large enough. Moreover  $k$ , the isomorphism class of the algebras, the trace and the inclusion matrices are all determined by the standard matrix of  $N \subset M$  and  $M \subset M_1$  ([Po2, Oc2]). An axiomatisation of such canonical commuting squares, in the finite depth case, describing them as paragroup type objects, was obtained in [Oc2]. In the case of the examples  $N_{G,\sigma} \subset M$  in 1.3.3 with  $\sigma$  an action of the group  $G$ , the matrix  $A$  is just the Cayley matrix of  $G$  and it does describe completely the canonical invariant of the inclusion ([Po4]). In particular, if  $G$  is infinite, then  $A$  is infinite and so  $N \subset M$  has infinite depth.

We mention that the basic construction can also be performed in a reverse way, this way obtaining a decreasing "tunnel" of subfactors  $M \supset N \supset N_1 \supset \dots$ . The existence of each subfactor  $N_{k+1}$  in this tunnel is however unique only up to conjugacy by a unitary in  $N_k$  ([PiPo1]). But the isomorphism class of the sequence of finite dimensional inclusions  $\{N'_k \cap N \subset N'_k \cap M\}_k$  as well as their resulting closures  $R_0 = \overline{\cup(N'_k \cap N)} \subset \overline{\cup(N'_k \cap M)} = R$  depends only on the class of  $N \subset M$  and is an obvious candidate for a model for  $N \subset M$  ([Po2]).

Classifying  $N \subset M$  by their canonical invariant amounts to proving the existence of a choice of a tunnel  $M \supset N \supset N_1 \supset \dots$  for which the higher relative commutants  $N'_k \cap M$  generate  $M$ . We prove such an exhaustion result for subfactors for which the higher relative commutants satisfy a certain "amenable growth" condition.

**2.2 Definition** [Po3].  $N \subset M$  is *strongly amenable* if the  $k$ 'th higher relative commutant for  $R_0 \subset R$  is isomorphic to the  $k$ 'th higher relative commutant for  $N \subset M$ , for each  $k \geq 1$ . This condition is also equivalent to the similar one

obtained by considering the inclusion  $M'_1 \cap M_\infty \subset M' \cap M_\infty$  instead of  $R_0 \subset R$ , where  $M_\infty = \overline{\cup M_k}$ . We mention that in fact  $\dim R'_0 \cap R = \dim N' \cap M$  (or equivalently  $\dim (M'_1 \cap M_\infty)' \cap M' \cap M_\infty = \dim N' \cap M$ ) is sufficient for the strong amenability condition to hold true.

**2.3 Theorem** [Po3]. *Assume  $N \subset M$  are hyperfinite type II<sub>1</sub> factors with finite index. The following conditions are equivalent:*

- 2.3.1  $N \subset M$  is strongly amenable.
- 2.3.2  $N \subset M$  is isomorphic to  $R_0 \subset R$ .
- 2.3.3 The bicommutant of  $M$  in  $M_\infty = \overline{\cup M_k}$  is equal to  $M$ .
- 2.3.4  $H(M \mid N) = H(R \mid R_0)$ , where, for  $A \subset B$ ,  $H(B \mid A)$  denotes the Connes-Størmer relative entropy [CS] as used in [PiPo1].

Moreover, the isomorphism class of a strongly amenable inclusion is completely determined by its canonical commuting squares.

Condition 2.3.4 in the above theorem can be interpreted as a Shannon-McMillan-Breiman type condition on the random walk with transition matrix  $A^t A$  ( $A$  being the standard matrix for  $N \subset M$ ) and transition probabilities proportional to the local indices of  $N \subset M_k$ : this random walk must have the largest possible entropy. The bicommutant property 2.3.3 is similar to the one stated for subfactors with finite depth and trivial relative commutant in [Oc2] and first pointed out in the case  $[M : N] < 4$  by Skau [GHJ].

The proof of the theorem uses much of the subfactor techniques developed in [PiPo1], especially the probabilistic characterization of the index  $[M : N]$  as the best constant  $s$  for which the inequality  $E_N(x) \geq s^{-1}x$  holds true for all  $x \in M_+$ ,  $E_N$  being the trace preserving conditional expectation of  $M$  onto  $N$ . The proof also uses in a crucial way the noncommutative local Rohlin theorem in [Po1].

The finite depth condition  $\dim A < \infty$  (equivalently  $\sup \dim \mathcal{Z}(M' \cap M_k) < \infty$ ,  $\mathcal{Z}(B)$  being the center of  $B$ ) is easily seen to imply the strong amenability condition. The particular case of Theorem 2.3 showing that subfactors with finite depth are classified by their canonical commuting squares was already settled in [Po2, Oc2].

All subfactors of index  $< 4$  have finite depth. They are thus classified by their canonical invariants. A full list of the combinatorial objects arising this way from subfactors of index  $< 4$  has been given in [Oc2] (for the case  $E_6$  see also [BN]).

Condition 2.3.4 implies that all subfactors of index 4 are strongly amenable, although some do not have finite depth. A full list of the combinatorial objects that classify them, and thus of all subfactors of index 4 is given in [Po3].

**2.4 Corollary.** *Subfactors of index  $\leq 4$  are all strongly amenable and are thus completely classified by their canonical commuting squares. A full list of these subfactors can be given.*

Another important class of subfactors to which 2.3 applies are the Wenzl subfactors 1.3.4, which have finite depth and are thus strongly amenable, and the Wassermann subfactors 1.3.5, which are strongly amenable but have infinite

depth. All these subfactors are thus uniquely determined by their combinatorial data.

### 3. Other Applications

Due to example 1.3.3 the classification Theorem 2.3 implies the classification, up to cocycle conjugacy, of actions  $\sigma$  of finitely generated discrete groups  $G$  on the hyperfinite factor, in the case when the associated subfactors  $N_{G,\sigma} \subset M$  result strongly amenable. In fact, 2.3 translates into :

**3.1 Theorem [Po4].** *Let  $P$  be the hyperfinite type  $\text{II}_1$  factor and  $\sigma$  an outer action of a finitely generated discrete group  $G$  on  $P$ . With the notations of 1.3.3 the following are equivalent.*

3.1.1  $N_{G,\sigma} \subset M$  is strongly amenable.

3.1.2 If  $g_i$  are some generators of  $G$  and  $a = \frac{1}{2n+1}(\sum g_i + e + \sum g_i^{-1})$  then  $\|g * a^n - a^n\|_1 \rightarrow 0$  for all  $g \in G$ .

3.1.3 The group  $G$  has 0 entropy, in the sense of [KV].

Moreover, the class of  $N_{G,\sigma} \subset M$  depends only on  $G$  and not on the cocycle conjugacy class of  $\sigma$ . Thus, there is a unique action, up to cocycle conjugacy, of the group  $G$  on the hyperfinite type  $\text{II}_1$  factor.

Note that the groups  $G$  satisfying 3.1.2 or 3.1.3 are automatically amenable. The above theorem thus implies a particular case from Ocneanu's theorem on the uniqueness of actions of amenable groups (in fact, with a slight modification, for  $\text{II}_\infty$  factors as well). The class of groups with 0 entropy contains all groups with subexponential growth and thus most classes of important amenable groups. In particular 3.1 (and thus the general Theorem 2.3) does cover Connes' results on the classification of single automorphisms, which is essential for the classification of hyperfinite von Neumann algebras.

The relation between the ergodic properties of the inclusion of  $N_G \subset M$  and of the random walk on  $G$  are investigated in [B].

Another application of 2.3 is the classification of minimal actions of compact groups.

**3.2 Theorem [PoWa].** *Given any compact Lie group  $G$  there exists a unique outer minimal action of  $G$  on the hyperfinite factor.*

To prove this result one shows that, up to conjugacy, a minimal action of  $G$  is uniquely determined by the isomorphism class of the inclusion of factors  $M^G \subset (M \otimes \text{End } V)^G$  of 1.3.5, with  $\pi$  a finite dimensional representation of  $G$  containing a generating set of irreducible representations. This inclusion follows strongly amenable by the Wassermann's invariance principle ([Wa]) and so 2.3 applies.

The subfactor techniques can also be used to give a simple elementary proof to an old standing problem on cocycle actions:

**3.3 Theorem** [Po4]. *Let  $G$  be a discrete group with 0 entropy (e.g. with subexponential growth). Then any 2-cocycle action of  $G$  on an arbitrary type II factor can be perturbed to a genuine action.*

We mention that one can prove Theorem 2.3 by using only the injectivity property of the ambient factor  $M$  (i.e. for  $M$  having a certain invariant mean, called hypertrace). If one takes  $N \subset M = M_{2 \times 2}(N)$  then such a general form of 2.3 implies the hyperfiniteness of the injective factor  $M$  and thus the uniqueness of the hyperfinite type  $\text{II}_{\infty}$  factor, the other major result of Connes that led to the classification of hyperfinite algebras.

One can furthermore obtain an equivariant version of Theorem 2.3, involving an automorphism  $\theta$  acting on  $M$  and leaving  $N$  globally invariant. The result shows that, aside from a commonly splitted part,  $\theta$  acts only on the combinatorial part of  $N \subset M$ , resulting from the higher relative commutants. Due to the work in ([L]) this implies the classification of strongly amenable subfactors of type  $\text{III}_{\lambda}$ ,  $0 < \lambda < 1$ , in terms of their combinatorial data.

## 4. Gaps for the Index of Hyperfinite Subfactors

Theorem 2.3 shows that not all hyperfinite type  $\text{II}_1$  subfactors can be approximated (and thus classified) by finite dimensional subalgebras coming from higher relative commutants. It is in fact likely that not all hyperfinite type  $\text{II}_1$  subfactors can be obtained by approximation with finite dimensional commuting squares like in example 1.3.6 (see 0.3 in [Po4]). For such subfactors though, one may hope that a classification can still be completed by establishing necessary and sufficient conditions for two increasing sequences of finite dimensional commuting squares (1.3.6) to give isomorphic subfactors, and then by classifying such sequences.

Approximation by finite dimensional commuting squares is also important for investigating the set of indices of irreducible subfactors, as subfactors with such approximation properties have index equal to square norms of matrices (cf. 1.3.6).

There is a rather general class of hyperfinite subfactors for which such approximations can be obtained. The key ingredient is the existence of certain special hypertraces  $\phi$  for  $N \subset M$ , called *symmetric hypertraces*, for which both  $M$ ,  $M'$  and the expectation of  $M$  onto  $N$ ,  $e_N$ , are contained in the centralizer of  $\phi$ . To prove existence of such hypertraces one needs  $N \subset M$  to have certain ergodicity properties, much weaker though than the Shannon-Mc Millan-Breiman condition in 2.3.4.

**4.1 Definition.**  $N \subset M$  is called *weakly amenable* if the random walk associated to it is ergodic, i.e. if the higher relative commutants generate a factor.

Note that all subfactors of sufficiently small indices ( $< 2 + \sqrt{5}$  will do) follow automatically weakly amenable.

**4.2 Theorem** [Po6]. *If  $N \subset M$  is a weakly amenable inclusion of hyperfinite type  $\text{II}_1$  factors and  $N' \cap M = \mathbb{C}$ , then  $N \subset M$  has a symmetric hypertrace.*

Using the techniques in ([Po7]) one then proves:

**4.3 Theorem** [Po6]. *If  $N \subset M$  is a weakly amenable inclusion of hyperfinite type  $\text{II}_1$  factors with  $N' \cap M = \mathbb{C}$ , then  $N \subset M$  can be locally approximated by finite dimensional commuting squares.*

By [PiPo2] (see 1.3.6) and taking into account that subfactors for which  $[M : N] < 2 + \sqrt{5}$  are all weakly amenable and that any number  $\geq 2 + \sqrt{5}$  is the square norm of an integral matrix ([GHJ]) one gets:

**4.4 Corollary** [Po6]. *If  $N \subset M$  are hyperfinite  $\text{II}_1$  factors with  $N' \cap M = \mathbb{C}$ , then  $[M : N]$  is equal to the square of the norm of a (possibly infinite) matrix with nonnegative integer entries.*

The set of square norms of (possibly infinite) matrices with nonnegative integer entries is known (see e.g. [GHJ]): besides containing the half line  $[2 + \sqrt{5}, \infty]$ , it has only countable accumulation points below  $2 + \sqrt{5}$ , which converge to  $2 + \sqrt{5}$ . So it has plenty of gaps below  $2 + \sqrt{5}$ . The first gap arises between 4 and the square norm of the matrix corresponding to the Coxeter graph  $E_{10}$ ,  $\|E_{10}\|^2 = 4.026\dots$ , as this is the first matrix with square norm larger than 4. We mention that the existence of irreducible hyperfinite subfactors of such index was proved in [HOcS].

**4.5 Corollary.** *The set of indices of irreducible subfactors of the hyperfinite type  $\text{II}_1$  factor has a gap between 4 and 4.026.*

In order to characterize completely the set of indices of irreducible subfactors of the hyperfinite  $\text{II}_1$  factor one should be able to construct examples as well. Some examples were obtained by constructing commuting squares of finite dimensional algebras, like in 1.3.6, in [GHJ, We1, We2, HOcS, Ch, Su].

We mention that the inclusion matrices for the finite dimensional approximations of the weakly amenable hyperfinite type  $\text{II}_1$  factors  $N \subset M$  can be interpreted more “canonically”, by introducing the *universal graph* (or *matrix*) of  $N \subset M$ , a general invariant that we will not explain here.

## 5. The Nonhyperfinite Case

Quite surprisingly, unlike the hyperfinite case when, as we have seen, the situation is quite rigid and the universal matrix of the subfactor is forced to have square norm equal to the index, for general irreducible pairs of type  $\text{II}_1$  factors any index  $s > 4$  may occur ([Po5]).

We will briefly discuss here a general method for constructing subfactors of finite index, which is quite different from the ones discussed before in 1.3 (e.g.

by producing finite dimensional commuting squares). The method consists in constructing Markov traces on certain universal algebras associated to an algebra  $Q$  and to the Jones' sequence of projections  $\{e_i\}_i$ ,  $\tau(e_i) = s^{-1}$ . In particular, for certain extremal such traces, this produces one parameter families of irreducible inclusions of factors  $N^s \subset M^s$  of index  $[M^s : N^s] = s$ ,  $s$  ranging over the whole set  $\{4 \cos^2 \frac{\pi}{n} \mid n \geq 3\} \cup [4, \infty)$ . For these extremal traces the resulting factors  $M^s$  are always non  $\Gamma$  (in the sense of [MvN]) and thus nonhyperfinite. The subfactors  $N^s \subset M^s$  are not only irreducible, but also their higher relative commutants  $(N^s)' \cap M_k^s$  are in some sense minimal, as they are generated by just  $e_1, e_2, \dots, e_k$ .

**5.1 Theorem** [Po5]. *Given any  $s > 4$  there are irreducible inclusions of nonhyperfinite factors of index  $s$ .*

The construction of  $N^s \subset M^s$  can be briefly described as follows: Let  $Q$  be an algebra. Consider the universal algebra  $U^s$  generated by  $Q$  and by the Jones' projections  $\{e_i\}_i$  of trace  $\tau(e_i) = s^{-1}$ , subject to the commutation relation  $[Q, e_i] = 0$ , for  $i \geq 2$ . Define on  $U^s$  a trace by letting  $tr(w) = 0$  for all words with alternating letters  $x_i \in Q$ ,  $y_j \in \text{Alg}\{1, e_1, e_2, \dots\}$ , for which  $\tau(x_i) = 0$  and for which the projection of  $y_j$  onto  $\text{Alg}\{1, e_2, \dots\}$  is 0 ([V]). Then  $M^s$  is defined as the smallest algebra of  $U^s/tr$  containing  $Q$  and on which  $e_1$  implements by compression a conditional expectation and  $N^s$  is defined as the commutant of  $e_1$  in  $M^s$ . The fact that  $[M^s : N^s] = s$  results from a certain remarkable property of the above trace, called the Markov property.

In fact, any pair of factors arises this way, from some  $Q$  and some appropriate Markov trace on  $U^s = U^s(Q)$ . The free Markov trace defined before correspond to an extremal situation when no “information” is exchanged from  $Q$  to  $\text{Alg}\{e_1, e_2, \dots\}$  through the “window”  $e_1$ . How to practically construct some others, especially hyperfinite ones, irreducible and with a prescribed index, is an open problem. This however may turn out to be a better method for constructing hyperfinite subfactors than the construction of finite dimensional commuting squares.

## References

- [B] D. Bisch: Entropy of groups and subfactors. To appear in Journal of Functional Analysis
- [BN] J. Bion-Nadal: Subfactors associated to  $E_6$ . CNRS, Preprint 1989
- [C] A. Connes: Classification des facteurs. Proc. Symp. Pure Math. **38** (1982) 43–109
- [Ch] M. Choda: Index for factors generated by Jones' two sided sequence of projections. Pac. J. Math. **139** (1989) 1–16
- [CS] A. Connes, E. Størmer: Entropy for automorphisms of  $\text{II}_1$  von Neumann algebras. Acta Math. **134** (1975) 288–259
- [GHJ] F. Goodman, P. de la Harpe, V.F.R. Jones: Coxeter graphs and tower of algebras. MSRI Publications 14. Springer 1989
- [HOcS] U. Haagerup, A. Ocneanu, J. Schou: In preparation
- [J1] V.F.R. Jones: Actions of finite groups on the hyperfinite  $\text{II}_1$  factor. Mem. A.M.S. **28**, no. 237 (1980)

- [J2] V.F.R. Jones: Index for subfactors. *Invent. math.* **72** (1983) 1–25
- [J3] V.F.R. Jones: Hecke algebra representations of braid groups and link polynomials. *Ann. Math.* **126** (1987) 335–388
- [JT] V.F.R. Jones, M. Takesaki: Actions of compact abelian groups on semifinite injective factors. *Acta math.* **153** (1984) 213–258
- [KV] V. Kaimanovich, A. Vershik: Random walks on discrete groups: boundary and entropy. *Ann. Probab.* **11** (1983) 457–490
- [L] P. Loi: On automorphisms of subfactors. UCLA, Preprint 1990
- [MvN] F. Murray, J. von Neumann: Rings of operators IV. *Ann. Math.* **44** (1943) 716–808
- [Oc1] A. Ocneanu: Actions of discrete amenable groups on von Neumann algebras. (Lecture Notes in Mathematics, vol. 1138.) Springer, Berlin Heidelberg New York 1985
- [Oc2] A. Ocneanu: Quantized groups string algebras and Galois theory for algebras. In: Operator Algebras and Applications, 2. London Math. Soc. Lect. Notes Series **136** (1988) 119–172
- [PiPo1] M. Pimsner, S. Popa: Entropy and index for subfactors. *Ann. Sci. Ec. Norm. Sup.* **19** (1986) 57–106
- [PiPo2] M. Pimsner, S. Popa: Finite dimensional approximation of pairs of algebras and obstructions for the index. To appear in *J. Funct. Anal.*
- [Po1] S. Popa: On a problem of R.V. Kadison on maximal abelian \*-subalgebras in factors. *Invent. math.* **65** (1981) 269–281
- [Po2] S. Popa: Classification of subfactors: reduction to commuting squares. *Invent. math.* **101** (1990) 19–43
- [Po3] S. Popa: Sur la classification des sous-facteurs d'indice fini du facteur hyperfini. *C. R. Acad. Sci. Paris, Série I* **311** (1990) 95–100
- [Po4] S. Popa: Sousfacteurs, actions des groupes et cohomologie. *C. R. Acad. Sci. Paris, Série I* **309** (1989) 771–776
- [Po5] S. Popa: Markov traces on universal Jones algebras and subfactors of finite index. IHES Preprint 1990
- [Po6] S. Popa: In preparation
- [Po7] S. Popa: On amenability in type  $\text{II}_1$  factors. In: Operator Algebras and Applications 2. London Math. Soc. Lect. Notes Series **136** (1988) 173–183
- [PoWa] S. Popa, A. Wassermann: In preparation
- [Su] S. Sunder: From hypergroups to subfactors. Preprint 1989
- [V] D. Voiculescu: Circular and semicircular systems and free product factors. UCB Preprint 1989
- [Wa] A. Wassermann: Coactions and Yang-Baxter equations for ergodic actions and subfactors. In: Operator Algebras and Applications 2. London Math. Soc. Lect. Notes Series **136** (1988) 202–236
- [We1] H. Wenzl: Hecke algebras of type  $A_n$  and subfactors. *Invent. math.* **92** (1988) 349–383
- [We2] H. Wenzl: Quantum groups and subfactors of type B, C and D. *Commun. Math. Phys.* **133** (1990) 383–432

# Operator Algebras and Duality

*Georges Skandalis*

Report on joint work with *Saad Baaj*

Université Paris 7, UFR de Math., URA 212, Tour 45-55, 2, place Jussieu  
F-75251 Paris Cedex 05, France

The Pontrjagyn duality associates to an abelian locally compact group a dual group and studies the properties of this correspondence. A natural idea is to try and generalize this duality to nonabelian groups, in particular to define an object dual to a group. Such dual objects were defined, first for compact groups [45, 19] then for locally compact groups ([38], [46]). Up to a large extent, it was the search of such a dual object for general locally compact groups that led to the theory of  $C^*$ -algebras. Then appeared the need of objects generalizing groups as well as their dual objects. These general objects can be called in a modern language “quantum groups”. These “groups” can be studied as abstract groups, Lie groups, deformations of true groups . . . . It is certainly beyond our goals to review all aspects of the theory of “quantum groups” (see [3] and references therein). We will in fact concentrate our attention to the operator algebra approach, in other words to the study of the “locally compact quantum groups”.

In terms of operator algebras Pontrjagyn duality takes the form of Takesaki-Takai duality [43, 40] based on the construction of  $W^*$ - and  $C^*$ -crossed-products. Along the years, under conjugated efforts of many specialists a set of axioms was built [11, 12, 42, 47, 16, 5] and duality was obtained [43, 20, 21, 30, 39, 4, 10, 6] for von Neumann algebras obeying these axioms called Kac von Neumann algebras. In a recent fundamental work [51–54], Woronowicz defined some objects that he called “compact matrix pseudogroups”. Although they aren’t Kac algebras, Woronowicz’ “pseudo-groups” enjoy duality properties. Further examples with the same properties were given by Majid [26] and Podles-Woronowicz [32]. One of the motivations of this report is to describe a setting including both the Kac von Neumann algebras and these new examples, in which the duality results still hold.

Let  $H$  be a Hilbert space and  $V \in L(H \otimes H)$  a unitary operator. Let us say  $V$  is multiplicative if it satisfies the pentagon equation  $V_{12}V_{13}V_{23} = V_{23}V_{12}$ . This relation appears in the framework of categories with associative tensor product (cf. [24, 25]); it is the one satisfied by the fusion operator (cf. [29]). It is also very similar to the Yang-Baxter equation and in some sense more primitive. In many papers concerned with operator algebras possessing duality properties, a multiplicative unitary plays a fundamental role (e.g. [12, 42, 16, 5, 18, 10] . . .); it is clear and more or less explicit in these papers that this unitary describes the whole situation.

It is therefore very natural to look for additional conditions that a multiplicative unitary should satisfy in order to correspond to a “locally compact quantum group”. Studying this problem, we were led to consider two conditions that we call *regularity* and *irreducibility*. We show how two pairwise dual Hopf  $C^*$ -algebras can be associated to a regular multiplicative unitary. When moreover this unitary is irreducible, we establish Takesaki-Takai duality results generalizing the previous ones.

An advantage of our approach is that a “quantum group” and its dual play completely symmetric roles. Also, it treats simultaneously the  $C^*$ - and  $W^*$ -algebra point of view. In fact, in our approach it is clear that for “locally compact quantum groups” the measure theory determines the topology. In some sense this can be thought of as a generalization of the famous theorem of Weil ([50], see also [23]): a “measurable quantum group” with an invariant (class of) measure(s) carries a unique structure of “locally compact quantum group”.

Let us also mention that many algebraic constructions can be performed in our setting. In particular, we may associate a “quantum double” to any (irreducible) multiplicative unitary, and together with it, comes a solution of the quantum Yang-Baxter equation (cf. [3]).

The question of the minimality of our axioms remains still unanswered: is a multiplicative unitary automatically regular? irreducible? Does one of these properties imply the other? Partial solutions to these questions were obtained: when the Hilbert space  $H$  is finite dimensional and when the unitary  $V$  satisfies a commutativity condition, regularity and irreducibility are both automatic; if the unitary  $V$  is of compact or discrete type (in other words if the associated quantum group is compact or discrete) its regularity implies its irreducibility.

In this report, we will first present some examples of occurrence of multiplicative unitaries, then explain the conditions of regularity and irreducibility and their consequences; we will finally construct the multiplicative unitaries associated with the examples of [26, 32] and discuss possible future developments. All the proofs, as well as more precise statements of the results given here can be found in [2].

## 1. Multiplicative Unitaries and Hopf Algebras

Let  $H$  be a Hilbert space. We will say that a unitary operator  $V$  acting on the tensor square  $H \otimes H$  is *multiplicative* if it satisfies the pentagon equation:

$$V_{23} V_{12} = V_{12} V_{13} V_{23}$$

Here, by  $V_{12}$ ,  $V_{23}$  and  $V_{13}$  we denote the operators  $V \otimes 1_H$ ,  $1_H \otimes V$  and  $(1_H \otimes \Sigma)(V \otimes 1_H)(1_H \otimes \Sigma)$  acting on  $H \otimes H \otimes H$ , where  $\Sigma$  is the “flip” operator defined by  $\Sigma(\xi \otimes \eta) = \eta \otimes \xi$  ( $\xi, \eta \in H$ ).

Note that the identity operator  $1_{H \otimes H}$  is a multiplicative unitary. The importance of multiplicative unitaries in connection with operator algebras possessing duality properties was shown by many authors [12, 42, 16, 5, 10]. The multiplicative unitary associated with a locally compact group is constructed as follows:

Let  $G$  be a locally compact group and let  $m$  denote its right Haar measure. Let  $H = L^2(G; m)$  be the Hilbert space of square integrable functions on  $G$  with respect with the measure  $m$ . Identify  $H \otimes H$  with the space  $L^2(G \times G; m \times m)$ . Let then  $V$  be the operator acting on  $H \otimes H$  by the formula  $V(f)(s, t) = f(st, t)$  for every square integrable function  $f$  on  $G \times G$  and  $s, t \in G$ . This operator is clearly unitary and its “multiplicativity” follows from the associativity of the composition law of  $G$ .

Operators satisfying this pentagon equation are naturally associated with Hopf algebras. Recall that a Hopf  $C^*$ -Algebra is a  $C^*$ -algebra  $A$  endowed with a coproduct which is a \*-homomorphism  $\delta : A \rightarrow A \otimes A^{(1)}$  satisfying the associativity condition:  $(\delta \otimes \text{id}) \circ \delta = (\text{id} \otimes \delta) \circ \delta$ . Let us describe three different ways for interpreting the pentagon relation:

### a) Haar States and GNS Representations

A Haar state on a Hopf  $C^*$ -algebra is a state  $\phi \in A^*$  such that for any form  $\psi \in A^*$  we have  $(\phi \otimes \psi) \circ \delta = (\psi \otimes \phi) \circ \delta = \psi(1)\phi$ . Let then  $(H_\phi, \pi_\phi, \xi_\phi)$  be the GNS construction associated with  $\phi$ . Then the operator  $V_\phi$  defined by  $V_\phi(\pi_\phi(x)\xi_\phi \otimes \eta) = (\pi_\phi \otimes \pi_\phi)(\delta(x))(\xi_\phi \otimes \eta)$  is an isometry of  $H_\phi \otimes H_\phi$  satisfying the pentagon equation. In particular, if  $V_\phi$  is surjective, it is a multiplicative unitary.

### b) Covariant Representations

Let  $(A, \delta)$  be a Hopf  $C^*$ -algebra. A corepresentation of  $A$  in a Hilbert space  $H$  is a unitary  $u \in L(H \otimes A)$  of the Hilbert  $A$ -module  $H \otimes A$  satisfying the relation:  $(\text{id} \otimes \delta)(u) = u_{12}u_{13}$ .

A coaction of  $A$  on some  $C^*$ -algebra  $B$  is a \*-homomorphism  $\delta_B : B \rightarrow B \otimes A$  satisfying the associativity condition:  $(\delta_B \otimes \text{id}) \circ \delta_B = (\text{id} \otimes \delta) \circ \delta_B$ . A covariant representation of  $A, B$  on a Hilbert space  $H$  is a pair  $(\pi, u)$  where  $\pi : B \rightarrow L(H)$  is a \*-representation and  $u \in L(H \otimes A)$  is a corepresentation of  $A$  such that  $\forall b \in B$ ,  $(\pi \otimes \text{id}) \circ \delta_B(b) = u(\pi(b) \otimes 1)u^*$ .

The coproduct  $\delta$  is a coaction of  $A$  on itself. Let  $(\pi, u)$  be a covariant representation on the Hilbert space  $H$ . Then  $V = (\text{id} \otimes \pi)(u)$  is a multiplicative unitary.

### c) The Canonical Element

Let  $A$  be a finite dimensional Hopf algebra. Let  $E$  be the algebra of endomorphisms of the vector space  $A$ . Let us denote by  $v$  the canonical element of  $A^* \otimes A$ : through the identification of  $A^* \otimes A$  with  $E$ ,  $v$  is the identity of  $A$ . Denote by  $L$  the action of  $A$  on  $A$  by left multiplication. If  $x \in A^*$  and  $a \in A$ , set  $\varrho(x)a = (\text{id} \otimes x)\delta(a)$ . Consider  $L$  and  $\varrho$  as homomorphisms from  $A$  and  $A^*$  into the algebra  $E$ . Simple computations then show:

---

<sup>1</sup> This is the  $C^*$ -algebraic “min” tensor product. If  $A$  has no unit,  $\delta$  takes its values in the multiplier algebra  $M(A \otimes A)$  of  $A \otimes A$ .

- For  $a \in A$  we have in  $E \otimes A$ ,  $(\varrho \otimes \text{id})(v)(L(a) \otimes 1) = (L \otimes \text{id})(\delta(a))(\varrho \otimes \text{id})(v)$ .
- We have the equality  $(\text{id} \otimes \delta)(v) = v_{12}v_{13}$ .

Therefore, the operator  $V = (\varrho \otimes L)(v)$  satisfies the pentagon equation.

## 2. Algebras Associated with Multiplicative Unitaries

Let  $V \in L(H \otimes H)$  be a multiplicative unitary. If  $\omega$  is a continuous linear form on  $L(H)$ , we may form the operators  $L(\omega) = (\omega \otimes \text{id})(V) \in L(H)$  and  $\varrho(\omega) = (\text{id} \otimes \omega)(V) \in L(H)$ . Let  $L(H)_*$  denote the predual of  $L(H)$  i.e. the set of linear mappings of the form  $x \rightarrow \text{Tr}(xT)$  where  $T$  spans the space of trace class operators.

**2.1. Proposition.** *The sets  $A(V) = \{L(\omega)/\omega \in L(H)_*\}$  and  $\hat{A}(V) = \{\varrho(\omega)/\omega \in L(H)_*\}$  are subalgebras of  $L(H)$ .*

Indeed,  $L(\omega)L(\omega') = (\omega \otimes \omega' \otimes \text{id})(V_{13}V_{23}) = L(\psi)$  where  $\psi(x) = (\omega \otimes \omega') \times (V^*(1 \otimes x)V)$

since  $V_{12}V_{13} = V_{23}V_{12}V_{23}^*$ ; in the same way,  $\varrho(\omega)\varrho(\omega') = (\text{id} \otimes \omega \otimes \omega')(V_{12}V_{13}) = \varrho(\psi')$  where  $\psi'(x) = (\omega \otimes \omega')(V(x \otimes 1)V^*)$ . In fact, all properties which may be proved for  $A(V)$  are automatically proved for  $\hat{A}(V)$  since  $\Sigma V^* \Sigma$  is a multiplicative unitary.

There is a natural duality between  $A(V)$  and  $\hat{A}(V)$  expressed by the equalities  $\langle L(\omega), \varrho(\omega') \rangle = \omega(\varrho(\omega')) = \omega'(L(\omega)) = (\omega \otimes \omega')(V)$ .

It is also natural to consider the norm closures of the algebras  $A(V)$  and  $\hat{A}(V)$  that we denote by  $S_V$  and  $\hat{S}_V$ . It is not clear whether these are always  $C^*$ -algebras ie. if they are closed under the involution  $x \rightarrow x^*$  of  $L(H)$ . For this reason, we are led to make, in the next sections, some extra assumptions.

In the case of the multiplicative unitary associated with a group, the algebras  $A(V)$  and  $\hat{A}(V)$  are respectively the Fourier algebra  $A(G)$  acting by multiplication on  $L^2(G)$  and  $L^1(G)$  acting by (right) convolution on  $L^2(G)$ . Also  $S_V$  is the abelian  $C^*$ -algebra  $C_0(G)$  of continuous functions vanishing at infinity and  $\hat{S}_V$  the reduced  $C^*$ -algebra of the group  $G$ . In particular the Gelfand spectrum of  $S_V$  is  $G$ : we already have recovered  $G$  out of the associated multiplicative unitary. In fact, we get a converse to this statement:

**2.2 Theorem.** *If the associated algebra  $A(V)$  is commutative, the multiplicative unitary  $V$  is (up to multiplicity) the multiplicative unitary associated with a locally compact group.*

This theorem is a generalization theorem of [50, 23, 12, 41, 42, 47, 5, 52]. Of course, this theorem also classifies the multiplicative unitaries for which  $\hat{A}(V)$  is commutative, since this is equivalent to saying that  $A(\Sigma V^* \Sigma)$  is commutative.

Let us mention another case where no extra assumptions are needed:

**2.3 Theorem.** *A multiplicative unitary acting on a finite dimensional Hilbert space is (up to multiplicity) the multiplicative unitary associated with a finite dimensional Kac von Neumann algebra.*

### 3. Regularity; the “Compact” Case

Let us begin with a rather easy fact:

**3.1 Proposition.** *The set  $\mathcal{C}(V) = \{(\text{id} \otimes \omega)(\Sigma\omega) / \omega \in L(H)_*\}$  is a subalgebra of  $L(H)$ .*

Studying this algebra in the case of locally compact groups and more generally in the examples to be discussed below, we find that this algebra is formed of compact operators and is norm dense in the algebra of compact operators. This leads to the following definition:

**3.2 Definition.** *We will say that the multiplicative unitary  $V$  is regular if the norm closure of  $\mathcal{C}(V)$  coincides with the algebra  $K(H)$  of compact operators of  $H$ .*

Regularity turns out to be extremely efficient in proving nice properties of the associated algebras:

**3.3 Theorem.** *Let  $V$  be a regular multiplicative unitary. Then the algebras  $S$  and  $\hat{S}$  are Hopf  $C^*$ -algebras with coproducts given by  $\delta(x) = V(x \otimes 1)V^*$  and  $\hat{\delta}(y) = V^*(1 \otimes y)V$  ( $x \in S, y \in \hat{S}$ ). The operator  $V$  is a multiplier of the (spatial) tensor product  $\hat{S} \otimes S$ .*

This last property means that the closed subalgebra of  $L(H \otimes H)$  generated by  $y \otimes x, x \in S, y \in \hat{S}$  is closed under left and right multiplication by  $V$ . It is quite natural and helpful. In particular, it allows us to consider  $S$  and  $\hat{S}$  as abstract  $C^*$ -algebras and still make sense of  $V$  in every representation.

We are also in position to form crossed products for algebras with coactions of the Hopf algebra  $S$ : if a  $C^*$ -algebra  $A$  is endowed with a coaction  $\delta_A : A \rightarrow A \otimes S$  of  $S$ , the (reduced) crossed product  $A \times \hat{S}$  is the  $C^*$ -algebra of operators acting on the Hilbert  $A$ -module  $A \otimes H$  generated by the products of the form  $\delta_A(a)(1 \otimes y)$ ,  $a \in A, y \in \hat{S}$ . For  $a \in A, y \in \hat{S}$ ,  $\delta_A(a)$  and  $(1 \otimes y)$  are multipliers of  $A \times \hat{S}$ . We thus get homomorphisms  $\pi$  and  $\hat{\theta}$  from  $A$  and  $\hat{S}$  respectively into the multiplier algebra of  $A \times \hat{S}$ . Still our set of axioms is not complete in order to allow us to prove the suitable duality. On the other hand, this duality may now be proved in the “compact” case.

**3.4 Definition.** *A multiplicative unitary is said to be of compact type if the unit operator belongs to the algebra  $A(V)$ .*

If  $V$  is a multiplicative unitary associated with a compact group or with a Haar state of a unital Hopf algebra, it is of compact type. In a recent fundamental work

[51–54], S.L. Woronowicz introduced a set of axioms for “quantum groups” generated by a finite dimensional unitary representation. Woronowicz initially called his objects “compact matrix pseudogroups” but they are referred to as “compact quantum Lie groups”. Nice and tractable examples were produced [51, 53]. It is natural to define “compact quantum groups” as projective limits (this corresponds to inductive limits for the Hopf algebra of functions) of “compact quantum Lie groups”. It was shown in [52] that “compact quantum groups” possess a Haar state. It is quite clear that the corresponding operator is unitary. Its regularity is also easy. The converse to these facts is true, namely:

**3.5 Theorem.** *A regular multiplicative unitary of compact type is (up to multiplicity) the multiplicative unitary associated with a “compact quantum group” of Woronowicz.*

#### 4. Irreducibility and Takesaki-Takai Duality

In order to introduce the last condition needed for the duality, let us examine again the case of locally compact groups: we have been able to produce out of the multiplicative unitary associated with a locally compact group, the multiplication operators and the right regular representation. On  $L^2(G)$  acts moreover the left regular representation; moreover, left and right regular representations are equivalent and intertwined by a unitary operator  $U$  given by  $(U\xi)(g) = \Delta(g)^{1/2}\xi(g^{-1})$ , where  $\Delta$  is the module of the group.

This leads us to assume the existence of an operator  $U$  satisfying some equations:

**4.1 Definition.** a) A multiplicative unitary  $V \in L(H \otimes H)$  is said to be irreducible if there exists a unitary  $U \in L(H)$  such that  $U^2 = 1_H$ ,  $(V(U \otimes 1)\Sigma)^3 = 1_{H \otimes H}$  and such that the unitary  $\hat{V} = \Sigma(U \otimes 1)V(U \otimes 1)\Sigma$  is multiplicative.

b) A Kac system is a triple  $(H, V, U)$  where  $H$  is a Hilbert space,  $V \in L(H \otimes H)$  is a multiplicative unitary and  $U \in L(H)$  satisfies the requirements of a); moreover, we require that  $V$  and  $\hat{V}$  be regular.

If  $(H, V, U)$  is a Kac system,  $(H, \hat{V}, U)$  is also a Kac system and  $\hat{\hat{V}} = (U \otimes U)V(U \otimes U)$ . Taking the dual again, we find  $\tilde{V} = \Sigma(1 \otimes U)V(1 \otimes U)\Sigma$ ; a fourth time will give us back  $V$ . This is the well noticed but still somewhat mysterious period 4 periodicity. (Note that as  $\hat{V} = (U \otimes U)V(U \otimes U)$  and  $\tilde{V} = (U \otimes U)\hat{V}(U \otimes U)$  they are regular multiplicative unitaries).

We now have two representations of  $S$  and  $\hat{S}$  in  $H$ : we will denote by  $L : S \rightarrow L(H)$  and  $\varrho : \hat{S} \rightarrow L(H)$  the inclusions considered up to now as identity representations; we will then set  $R(x) = UL(x)U$  and  $\lambda(y) = U\varrho(y)U$  ( $x \in S$ ,  $y \in \hat{S}$ ).

Replacing  $V$  by  $\hat{V}$  we may now form crossed products for algebras with coactions of the Hopf algebra  $\hat{S}$ : if a  $C^*$ -algebra  $A$  is endowed with a coaction  $\delta_A : A \rightarrow A \otimes \hat{S}$  of  $\hat{S}$ , the (reduced) crossed product  $A \times S$  is the  $C^*$ -algebra of operators acting on the Hilbert  $A$ -module  $A \otimes H$  generated by the products of the form

$$(\text{id} \otimes \lambda)\delta_A(a)(1 \otimes L(x)), \quad a \in A, \quad x \in S.$$

We still have homomorphisms  $\pi$  and  $\theta$  from  $A$  and  $S$  respectively into the multiplier algebra of  $A \times S$ , and  $A \times S$  is spanned by products  $\pi(a)\theta(x)$   $a \in A$ ,  $x \in S$ .

Let the  $C^*$ -algebra  $A$  be endowed with a coaction  $\delta_A$  of  $S$  (resp.  $\hat{S}$ ), then the crossed product  $A \times \hat{S}$  (resp.  $A \times S$ ) is endowed with a coaction  $\hat{\delta}_A$  of  $\hat{S}$  (resp.  $S$ ) given by  $\hat{\delta}_A(\pi(a)\hat{\theta}(y)) = (\pi(a) \otimes 1)(\hat{\theta} \otimes \text{id})\hat{\delta}(y)$ ,  $a \in A$ ,  $y \in \hat{S}$  (resp.  $\hat{\delta}_A(\pi(a)\theta(x)) = (\pi(a) \otimes 1)(\theta \otimes \text{id})\delta(x)$ ,  $a \in A$ ,  $x \in S$ ).

**4.2 Theorem** (Takesaki-Takai Duality Theorem). *Let  $(H, V, U)$  be a Kac system. Then for any algebra  $A$  endowed with a coaction<sup>2</sup>  $\delta_A$  of  $S$  the double crossed product  $A \times \hat{S} \times S$  is naturally isomorphic with  $A \otimes K(H)$ .*

*Remark.* Replacing  $(H, V, U)$  by  $(H, \hat{V}, U)$  we may exchange the roles of  $S$  and  $\hat{S}$ .

It is transparent in many papers (cf. eg. [21], [22], [39], [4], [10]) that Takesaki duality only relies on the “fundamental” operator; our proof is therefore just an adaption of methods used by these authors.

A first step to this duality is the case  $A = \mathbf{C}$ :

**4.3 Lemma.** *Let  $(H, V, U)$  be a Kac system. Then the closed vector span of  $\{L(x)\varrho(y)/x \in S, y \in \hat{S}\}$  and the closed vector span of  $\{L(x)\lambda(y)/x \in S, y \in \hat{S}\}$  are the algebra  $K(H)$ .*

This lemma, which can be thought of as a generalization of the famous Weyl-von Neumann theorem, explains the terminology of irreducibility.

*Remark.* It is also quite easy to prove a Takesaki duality theorem in the von Neumann algebra setting for Kac systems. In fact the regularity can be replaced by the weaker condition: the weak closure of  $\mathcal{C}(V)$  is  $L(H)$ . The proof (if not the precise statement) of the main theorem of [10] applies in this context.

Also, one may generalize results of [1] and prove a Takesaki-Takai duality theorem for equivariant  $KK$ -theory with respect to the Hopf  $C^*$ -algebras  $S$  and  $\hat{S}$ .

## 5. Examples of Majid and Podles-Woronowicz

In [26] and [32] appeared a series of new constructions of interesting “quantum groups”. These “quantum groups” are not in general Kac von Neumann algebras but they can still be expressed by a multiplicative unitary; in this way Takesaki-Takai duality is just an easy check.

The algebraic setting in the examples of [26, 27] and [32] is that of matched pairs of Hopf algebras (cf. [23, 28]). To such a matched pair are associated two new Hopf algebras: the one “generated by the matched pair” and the “bicrossproduct”. These constructions were given in [37] and [28] in purely algebraic terms but may be performed in the multiplicative unitary setting. Examples of such matched pairs

---

<sup>2</sup> Provided a technical assumption called non-degeneracy in [21] is fulfilled by  $\delta_A$ .

are given by any (locally compact) group  $G$  with two (closed) subgroups  $H$  and  $K$  such that every element  $g$  of  $G$  admits a unique decomposition  $g = hk$  ( $h \in H, k \in K$  – [44, 26]). Other examples are given by the “quantum double” construction of Drinfeld [3]. The examples of Podles and Woronowicz [32] are in fact based on this “quantum double” construction.

Let  $(A, \delta_A)$  and  $(B, \delta_B)$  be two Hopf  $C^*$ -algebras. Consider the \*-homomorphism  $\delta_A \otimes \delta_B : A \otimes B \rightarrow A \otimes A \otimes B \otimes B$ . In order to put a Hopf  $C^*$ -algebra structure on  $A \otimes B$ , we use a \*-isomorphism  $\tau : A \otimes B \rightarrow B \otimes A$  and wish to put  $\delta = (\text{id}_A \otimes \tau \otimes \text{id}_B)(\delta_A \otimes \delta_B)$ . For  $\delta$  to be coassociative it is enough that the following condition be satisfied:

$$(C) (\tau \otimes \text{id}_A)(\text{id}_A \otimes \tau)(\delta_A \otimes \text{id}_B) = (\text{id}_B \otimes \delta_A)\tau \quad \text{and} \\ (\text{id}_B \otimes \tau)(\tau \otimes \text{id}_B)(\text{id}_A \otimes \delta_B) = (\delta_B \otimes \text{id}_A)\tau$$

Condition (C) is stated in [32] and, from a dual point of view, in [37] and [28] (in purely algebraic terms).

**5.1 Definition.** Let  $(A, \delta_A)$  and  $(B, \delta_B)$  be two Hopf  $C^*$ -algebras. An inversion on  $A$ ,  $B$  is a \*-isomorphism  $\tau : A \otimes B \rightarrow B \otimes A$  satisfying the conditions (C).

Let  $(A, B, \tau)$  be as in definition 5.1. Note that  $b \rightarrow \tau(1 \otimes b)$  is a (right) coaction, called  $\beta$ , of the Hopf  $C^*$ -algebra  $A$  on the  $C^*$ -algebra  $B$  and  $a \rightarrow \tau(a \otimes 1)$  is a (left) coaction, called  $\alpha$ , of the Hopf  $C^*$ -algebra  $B$  on the  $C^*$ -algebra  $A$ .

Let  $X \in L(H \otimes H)$  and  $Y \in L(K \otimes K)$  be two regular multiplicative unitaries. Denote by  $S_X, \hat{S}_X, S_Y$  and  $\hat{S}_Y$  the associated Hopf  $C^*$ -algebras associated with  $X$  and  $Y$  and by  $\delta_X, \delta_Y$  the coproducts of  $S_X$  and  $\hat{S}_Y$ . Let  $\tau : S_X \otimes \hat{S}_Y \rightarrow \hat{S}_Y \otimes S_X$  be an inversion on  $(S_X, \hat{S}_Y)$ .

**5.2 Proposition.** The unitary operator  $T = (\tau \otimes \text{id})(Y_{23})(\text{id} \otimes \tau)(X_{23})$  acting on  $K \otimes H \otimes K \otimes H$  is multiplicative. It is called the bicrossproduct of  $X$  and  $Y$  with respect to  $\tau$ .

Let  $(H, X, u)$  and  $(K, Y, v)$  be two Kac systems. It is more natural to assume that  $\tau$  is an inversion on  $(S_X, S_Y)$ . Of course, in this case, we may form the bicrossproduct of  $X$  and  $\hat{Y}$ . In order to form the twisted and bicrossproducts of the Kac systems, we need the inversion  $\tau$  to be suitably implemented.

**5.3 Definition.** A matched pair of Kac systems is given by two Kac systems  $(H, X, u)$  and  $(K, Y, v)$  together with a unitary operator  $Z \in L(H \otimes K)$  such that

- a) There exists an inversion  $\tau : S_X \otimes S_Y \rightarrow S_Y \otimes S_X$  such that for all  $x \in S_X, y \in S_Y$  we have  $Z\tau^{-1}(y \otimes x)Z^* = x \otimes y$ .
- b)  $(H \otimes K, V, U)$  is a Kac system where  $V = (Z_{12}^* X_{13} Z_{12}) Y_{24}$  and  $U = (u \otimes v)Z$ .

**5.4 Theorem.** Let  $((H, X, u); (K, Y, v); Z)$  be a matched pair of Kac systems. Define the unitary operator  $W = (Z_{34}^* \hat{Y}_{24} Z_{34})(Z_{12}^* X_{13} Z_{12})$  acting on  $H \otimes K \otimes H \otimes K$ .

Then  $(H \otimes K, W, U)$  is a Kac system. The Hopf algebra  $S$  associated with  $V$  is  $(S_X \otimes S_Y, \delta_\tau)$ . The algebras  $S$  and  $\hat{S}$  associated with  $W$  are isomorphic respectively to  $S_X \times_\alpha S_Y$  and  $S_Y \times_\beta S_X$ .

Moreover, the multiplicative unitary operator  $T$  of proposition 5.2 is equivalent to  $W$ .

**5.5 Definition.** With the notation of the above theorem, the Kac system  $(H \otimes K, V, U)$  is called the product of  $(H, X, u)$  and  $(K, Y, v)$  twisted by  $Z$ ; the Kac system  $(H \otimes K, W, U)$  is called the bicrossproduct of  $(H, X, u)$  and  $(K, Y, v)$  relative to  $Z$ .

## 5.6 Examples

a) Let  $(H, X, u)$  be a Kac system and  $G$  be a locally compact group acting by Hopf  $C^*$ -algebra automorphisms on  $S_X$ . Let  $\tau : C_0(G) \otimes S_X \rightarrow S_X \otimes C_0(G)$  be given by  $\tau(f)(x) = \alpha_x(f(x))$   $x \in G$ ,  $f \in C_0(G; S_X)$  where we have identified  $C_0(G) \otimes S_X$  and  $S_X \otimes C_0(G)$  with the  $C^*$ -algebra  $C_0(G; S_X)$  of continuous  $S_X$ -valued functions vanishing at  $\infty$  on  $G$ . In this case,  $G$  acts naturally on the Hopf  $C^*$ -algebra  $\hat{S}_X$  and the twisted and bicrossproducts are both obtained by the well known crossed-product constructions.

b) Let  $G_1$  and  $G_2$  be two locally compact groups. An inversion on  $(C_0(G_1), C_0(G_2))$  is given by a homeomorphism  $\tau : G_2 \times G_1 \rightarrow G_1 \times G_2$ ; then the product  $(x_1, x_2)(y_1, y_2) = (x_1 z_1, z_2 x_2)$  where  $(z_1, z_2) = \tau(x_2, y_1)$  is associative on  $G_1 \times G_2$  and, endowed with this product,  $G_1 \times G_2$  is a locally compact group  $G$ . Then, the twisted product of the associated Kac systems is the Kac system of the group  $G$ . The bicrossproduct construction gives new examples of Kac systems. In general, these examples are not associated with Kac von Neumann algebras [26, 27] and the antipode  $\kappa$  is unbounded. However, many computations may still be performed in this context.

Another way of understanding this example, is to start with a locally compact group  $G$  and assume that it has two closed subgroups  $G_1$  and  $G_2$  such that the map  $(x_1, x_2) \rightarrow x_1 x_2$  is a homeomorphism from  $G_1 \times G_2$  onto  $G$ . In this case, the actions of  $G_1$  on  $G_2$  and of  $G_2$  on  $G_1$  are the restrictions of the actions of  $G$  on  $G_2 = G_1 \setminus G$  and on  $G_1 = G/G_2$  and it is easy to compute the corresponding crossed products and thus the algebras  $S$  and  $\hat{S}$  associated with  $W$ . Also, it is quite easy to construct groups with these properties:

- the Iwasawa decomposition  $G = KP$  ( $P = AN$ ) of semisimple Lie groups;
- let  $G$  be a locally compact group acting by homeomorphisms on a locally compact group  $G_2$  and containing the right translations of  $G_2$ ; let then  $G_1$  be the set of elements of  $G$  fixing the neutral element of  $G_2$ ;
- in the above example, we may take  $G_2$  to be any finite group and  $G$  be the group of all permutations of the set  $G_2 \dots$

A third way of interpreting this example (cf. [2] Appendix C) is the search of measure spaces  $X$  and transformations of  $X \times X$  satisfying the pentagon relation.

It is then natural to add cocycles and form new multiplicative unitaries. In this way one recovers examples of Kac and Paliutkin [14, 15].

c) Let  $(H, X, u)$  be a Kac system and set  $Y = sX^*s$  where  $s \in L(H \otimes H)$  is the flip operator; then  $(H, Y, u)$  is a Kac system and  $S_Y = \hat{S}_X$ ,  $\hat{S}_Y = S_X^{(3)}$ . Let then  $\tau : S_X \otimes S_Y \rightarrow S_Y \otimes S_X$  be given by  $\tau(x) = XsxX^*$ . Since  $X$  is a multiplier of  $S_Y \otimes S_X$  this is well defined. It turns out that  $\tau$  is a non degenerate inversion and that  $((H, X, u); (K, Y, v); Z)$  is a matched pair of Kac systems, where  $Z = sX(u \otimes u)X^*(u \otimes u)s$ . The corresponding twisted product  $(H \otimes H, V, U)$  is the *quantum double* of the Kac system  $(H, X, u)$ . Let  $S_V$  and  $\hat{S}_V$  denote the corresponding Hopf  $C^*$ -algebras. There is a unitary operator  $R$  which is a multiplier of  $\hat{S}_V \otimes \hat{S}_V$  which satisfies the algebraic properties of [3] and in particular  $R$  is a solution of the quantum Yang-Baxter equation.

Note that the construction of this twisted product was used by Podles and Woronowicz to build the “quantum  $SL(2, C)$ ” out of the “quantum  $SU(2)$ ” ([32]).

In this case, the bicrossproduct is just a direct product.

## 6. Concluding Remarks

We developed here one point of view: find conditions easy to check on the “fundamental operator” that ensure Takesaki-Takai duality. However, we do not know if these conditions may turn out to be automatic.

Maybe one should look for a counterexample to regularity in transformations satisfying the pentagon equation.

The operator  $U$  defining irreducibility, is usually the product  $\hat{J}J$  of the Tomita operators associated with Haar measures. Therefore, to prove irreducibility one would need to prove the existence of these Haar measures. Note that, in our context, this problem doesn’t seem too difficult since we are given the regular representations and therefore the class of the Haar measures.

Once the Haar measures are found one needs to perform modular theory on them. Concerning this, we may formulate the following conjecture:

Call  $\phi, \psi, \hat{\phi}$  and  $\hat{\psi}$  the left and right Haar measures of  $S$  and  $\hat{S}$ . Then there should exist positive unbounded operators  $F$  and  $\hat{F}$  affiliated with the centralizers of  $\phi$  and  $\hat{\phi}$  such that for all  $x \in S$  and  $y \in \hat{S}$ ,  $\psi(x) = \phi(Fx\bar{F})\hat{\psi}(y) = \hat{\phi}(\hat{F}\bar{y}\hat{F})$ . The Hilbert spaces  $H_\phi, H_\psi, H_{\hat{\phi}}$  and  $H_{\hat{\psi}}$  are naturally identified. The weights  $\phi, \psi, \hat{\phi}$  and  $\hat{\psi}$  are faithful when extended to the bicommutants, therefore Tomita theory can be performed. Call  $J$  and  $\hat{J}$  the Tomita operators associated with  $\phi$  and  $\hat{\phi}$ , and put  $U = \hat{J}J = J\hat{J}$ . Let  $L$  and  $\lambda$  be the GNS representations associated with  $\phi$  and  $\hat{\phi}$ ; then form  $R$  and  $\varrho$  using  $L, \lambda$  and  $U$ . Since  $V$  is a multiplier of  $\hat{S} \otimes S$ ,  $(\varrho \otimes L)(V)$  acts naturally on  $H_\phi \otimes H_\psi$ . Then  $(H_\phi, V, U)$  is a Kac system. Moreover, the operators  $F$  and  $\hat{F}$  are representations of the Hopf algebras, therefore they are unbounded multipliers of  $S$  and  $\hat{S}$ . Moreover, the operators  $L(F), R(F), \lambda(\hat{F})$  and  $\varrho(\hat{F})$  commute pairwise. The modular operators are computed in terms of  $F$  and  $\hat{F}$ . We find:

---

<sup>3</sup> Note however that the coproducts of  $S_Y$  and  $\hat{S}_Y$  differ from the ones of  $\hat{S}_X$  and  $S_X$  by the flip.



16. Kac, G.I., Vainerman, L.I.: Nonunimodular ring-groups and Hopf-von Neumann algebras. *Math. USSR Sb.* **23** (1974) 185–214. Translated from *Mat. Sb.* **94** (136) (1974), 2 194–225
17. Katayama, Y.: Takesaki's duality for a non degenerate coaction. *Math. Scand.* **55** (1985) 141–151
18. Kirchberg, E.: Representation of coinvolutive Hopf- $W^*$ -algebras and non abelian duality. *Bull. Acad. Pol. Sc.* **25** (1977) 117–122
19. Krein, M.G.: Hermitian-positive kernels in homogeneous spaces. *Amer. Math. Soc. Transl.* (2) **34** (1963) 109–164. Translated from *Ukr. Mat. Z.* **2**, no. 1 (1950) 10–59
20. Landstad, M.B.: Duality theory for covariant systems. *Trans. AMS* **248** (1979) 223–267
21. Landstad, M.B.: Duality for dual covariance algebras. *Comm. Math. Phys.* **52** (1977) 191–202
22. Landstad, M.B., Phillips, J., Raeburn, I., Sutherland, C.E.: Representations of crossed products by coactions and principal bundles. *Trans. AMS* **299**, no. 2 (1987) 747–784
23. Mackey, G.W.: Borel structures in groups and their duals. *Trans. AMS* **85** (1957) 134–165
24. Mac Lane, S.: Categories for the working mathematician. (Graduate Texts in Mathematics, vol. 5). Springer, New York Berlin Heidelberg 1974
25. Mac Lane, S.: Natural associativity and commutativity. *Rice Univ. Studies* **49** (1963) 4–28
26. Majid, S.H.: Non-commutative geometric groups by a bicrossproduct construction: Hopf algebras at the Planck scale. Thesis, Harvard Univ. 1988
27. Majid, S.H.: Hopf von Neumann algebra Bicrossproducts, Kac Algebras Bicrossproducts, and the classical Yang-Baxter equation. To appear in *J. Funct. Anal.*
28. Majid, S.H.: Physics for algebraists: Non-commutative and Non-cocommutative Hopf algebras by a Bicrossproduct construction. *J. Algebra* **130**, no. 1 (1990) 17–64
29. Moore, G., Seiberg, N.: Classical and quantum conformal field theory. *Comm. Math. Phys.* **123** (1989) 177–254
30. Nakagami, Y.: Dual action on a von Neumann algebra and Takesaki's duality for a locally compact group. *Publ. RIMS Kyoto Univ.* **12** (1977) 727–775
31. Nakagami, Y., Takesaki, M.: Duality for crossed products of von Neumann algebras. (Lecture Notes in Mathematics, vol. 731.) Springer, Berlin Heidelberg New York 1979
32. Podles, P., Woronowicz, S.L.: Quantum deformation of Lorentz group. Preprint
33. Rosso, M.: Comparaison des groupes  $SU(2)$  quantiques de Drinfeld et Woronowicz. *Note aux C. R. Acad. Sci.* **304** (1987) 323–326
34. Rosso, M.: Algèbres enveloppantes quantifiées, groupes quantiques compacts de matrices et calcul différentiel non commutatif. Prépublication
35. Schwartz, J.M.: Sur la structure des algèbres de Kac I. *J. Funct. Anal.* **34** (1979) 370–406
36. Schwartz, J.M.: Sur la structure des algèbres de Kac II. *Proc. London Math. Soc.* **41** (1980) 465–480
37. Singer, W.M.: Extension theory for connected Hopf algebras. *J. Algebra* **21** (1972) 1–16
38. Stinespring, W.F.: Integration theorems for gages and duality for unimodular groups. *Trans. AMS* **90** (1959) 15–56
39. Stratila, S., Voiculescu, D. and Zsidó, L.: On crossed products. I and II. *Rev. Roumaine Math. Pures Appl.* **21** (1976) 1411–1449 and **22** (1977) 83–117
40. Takai, H.: On a duality for crossed products of  $C^*$ -algebras. *J. Funct. Anal.* **19** (1975) 25–39
41. Takesaki, M.: A characterization of group algebras as a converse of Tannaka-Stinespring-Tatsuuma duality theorem. *Amer. J. Math.* **91** (1969) 529–564
42. Takesaki, M.: Duality and von Neumann algebras. (Lecture Notes in Mathematics, vol. 427.) Springer, Berlin Heidelberg New York 1972, pp. 665–779
43. Takesaki, M.: Duality for crossed products and the structure of von Neumann algebras of type III. *Acta Math.* **131** (1973) 249–310

44. Takeuchi, M.: Matched pairs of groups and bismash products of Hopf algebras. *Comm. Algebra* **9**, no. 8 (1981) 841–882
45. Tannaka, T.: Über den Dualität der nicht-kommutativen topologischen Gruppen. *Tôhoku Math. J.* **45** (1938) 1–12
46. Tatsuuma, N.: A duality theorem for locally compact groups. *J. Math. of Kyoto Univ.* **6** (1967) 187–293
47. Vainerman, L.I.: Characterization of dual objects for locally compact groups. *Funct. Anal. Appl.* **8** (1974) 66–67. Translated from *Funkt. Anal. Prilozhen.* **8**, no. 1 (1974) 75–76
48. Vallin, J.M.:  $C^*$ -algèbres de Hopf et  $C^*$ -algèbres de Kac. *Proc. London Math. Soc.* (3) **50** (1985) 131–174
49. Voiculescu, D.: Amenability and Katz algebras. *Algèbres d'opérateurs et leurs applications en physique mathématique. Colloques Internationaux C.N.R.S.* no. **274** (1977) 451–457
50. Weil, A.: L'intégration dans les groupes topologiques et ses applications. *Act. Sc. Ind.* no. **1145**. Hermann, Paris (1953)
51. Woronowicz, S.L.: Twisted SU(2) group. An example of a non commutative differential calculus. *Publ. RIMS* **23** (1987) 117–181
52. Woronowicz, S.L.: Compact matrix pseudogroups. *Comm. Math. Phys.* **111** (1987) 613–665
53. Woronowicz, S.L.: Tannaka-Krein duality for compact matrix pseudogroups. Twisted SU( $N$ ) group. *Inv. math.* **93** (1988)
54. Woronowicz, S.L.: Differential calculus on compact matrix pseudogroups (quantum groups). *Comm. Math. Phys.* **122** (1989) 125–170



# Some Isoperimetric Inequalities and Their Applications

*Michel Talagrand*

Equipe d'Analyse, E.R.A. au C.N.R.S. no. 294, Université Paris VI, 4 Place Jussieu  
F-75230 Paris Cedex 05, France

## 0. Introduction

Consider a product of measure spaces, provided with the product measure. Consider a subset  $A$  of this product, of measure at least one half. An important fact (the so-called concentration of measure phenomenon) is that even a small “enlargement” of  $A$  has measure very close to one. The inequalities we present describe this phenomenon for several notions of “enlargement”.

## 1. The Isoperimetric Inequality for Gaussian Measure

We denote by  $S_n$  the Euclidean sphere of  $\mathbb{R}^{n+1}$ , equipped with the geodesic distance  $\varrho$  and a rotation invariant probability  $\mu_n$ . For a (measurable) subset  $A$  of  $S_n$ , consider the set  $A_u$  of points of  $S_n$  within geodesic distance  $u$  of  $A$ . The isoperimetric inequality for the sphere, discovered by P. Lévy, is of fundamental importance. It states that  $\mu_n(A_u) \geq \mu_n(C_u)$ , where  $C$  is a cap (intersection of the sphere and of a half space) of the same measure as  $A$ .

We denote by  $\gamma_n$  the canonical Gaussian measure on  $\mathbb{R}^n$ , of density  $(2\pi)^{-n/2} e^{-\|x\|^2/2}$  with respect to Lebesgue measure. Observe the simple, but essential fact that  $\gamma_n$  is the product measure on  $\mathbb{R}^n$  when each factor is endowed with  $\gamma_1$ . It is an old observation, known as Poincaré lemma (although it does seem to be due to Maxwell) that, as  $N \rightarrow \infty$ , the projection of the normalized measure on  $\sqrt{N}S_N$  onto  $\mathbb{R}^n$  has  $\gamma_n$  as a limit. Therefore, it should not come as a surprise that Lévy's isoperimetric inequality on the sphere implies an isoperimetric inequality for  $\gamma_n$ . This was discovered independently by C. Borell [B1], and V. N. Sudakov and B. S. Tsirelson [S-T]. If we denote by  $A_u$  the set of points within Euclidean distance  $u$  of  $A$ , then  $\gamma_n(A_u) \geq \gamma_n(H_u)$ , where  $H_u$  is a half space with  $\gamma_n(H) = \gamma_n(A)$ . Taking this half space to be orthogonal to a coordinate axis, and remembering that  $\gamma_n$  is a product measure shows that if  $\gamma_n(H) = \gamma_1((-\infty, a])$ , then  $\gamma_n(H_u) = \gamma_1((-\infty, a + u])$ .

---

The author is also affiliated with the Ohio State University, Columbus, OH 43201, USA

For simplicity, we set  $\Phi(u) = \gamma_1((-\infty, u]) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u e^{-x^2/2} dx$ . Thus we have

$$\text{if } \gamma_n(A) = \Phi(a), \text{ then } \gamma_n(A_u) \geq \Phi(a+u). \quad (1.1)$$

It is important to state this inequality, not only for the measure  $\gamma_n$ , but also for its infinite dimensional version  $\gamma$ , the product measure on  $\mathbb{R}^{\mathbb{N}}$  when each factor is endowed with  $\gamma_1$  (the result for  $\gamma$  follows from the result for  $\gamma_n$  and an obvious approximation). We denote by  $B$  the unit ball of  $\ell^2$ , i.e.  $B = \{x \in \mathbb{R}^{\mathbb{N}}, \sum_{k \geq 1} x_k^2 \leq 1\}$ . The Gaussian isoperimetric inequality can then be stated as follows

$$\text{If } \gamma(A) = \Phi(a) \text{ then } \gamma_*(A + uB) \geq \Phi(a+u). \quad (1.2)$$

There  $A + uB = \{x + uy; x \in A, y \in B\}$ ; the inner measure is needed as  $A + uB$  might not be measurable. As became customary, we call (1.1) and (1.2) Borell's inequality. Lévy's inequality is usually proved using symmetrization (see e.g. the appendix of [F-L-M]). A. Ehrhard [E1] has developed a symmetrization method adapted to the measures  $\gamma_n$  that yields a direct proof of (1.2) as well as of the following remarkable inequality of Brunn-Minkowski's type: For two convex sets  $A, B$  of  $\mathbb{R}^n$ , and  $0 \leq \lambda \leq 1$ ,

$$\Phi^{-1}(\gamma_n(\lambda A + (1 - \lambda)B)) \geq \lambda \Phi^{-1}(\gamma_n(A)) + (1 - \lambda) \Phi^{-1}(\gamma_n(B)). \quad (1.3)$$

(It is still open whether this inequality holds for non convex sets.)

Borell's inequality is a principle of remarkable power. It can be argued that, concerning applications, this inequality is used in two different forms.

The first type of use consist of rewriting (1.1) as  $u^{-1}\gamma_n(A_u \setminus A) \geq u^{-1}\gamma_1([a, u+a])$  so that

$$\liminf_{u \rightarrow 0} u^{-1}\gamma_n(A_u \setminus A) \geq \frac{1}{\sqrt{2\pi}} \exp -\frac{a^2}{2}. \quad (1.4)$$

thereby recovering what is the more classical formulation of the isoperimetric inequality [O]. In this spirit (and using his symmetrization methods) A. Ehrhard has proved a number of interesting inequalities, that are versions for the Gauss measure of classical results [E2].

Inequality (1.4) for functions rather than sets [L] yields in particular that a function on  $\mathbb{R}^n$  whose gradient belongs to  $L^1(\gamma_n)$ , belongs to the Orlicz space  $L^1(\log L)^{1/2}$  of this measure, connecting with logarithmic Sobolev inequalities and hypercontractivity.

The second type of use of Borell's inequality is for “large” values of  $u$  (while Borell's inequality for large values of  $u$  follows from (1.4), the spirit of application is very different). It is mostly used in the following forms

$$\text{If } \gamma_n(A) \geq 1/2, \text{ then } \gamma_n(A_u) \geq \gamma_1((-\infty, u]) \quad (1.5)$$

$$\text{If } \gamma(A) \geq 1/2, \text{ then } \gamma_*(A + uB) \geq \gamma_1((-\infty, u]) \geq 1 - \frac{1}{2} \exp(-u^2/2). \quad (1.6)$$

In the terminology of V. Milman [M2] (1.5) is a “concentration of measure phenomenon”. An immediate consequence of (1.5) is that if  $f$  is a Lipschitz function on  $\mathbb{R}^n$ , we have

$$\gamma_n(\{|f - M_f| \geq u\}) \leq 2\gamma_1\left(\left[\frac{u}{\|f\|_{\text{Lip}}}, \infty\right)\right) \leq \exp\left(-\frac{u^2}{2\|f\|_{\text{Lip}}^2}\right) \quad (1.7)$$

where  $M_f$  is a median of  $f$ , i.e.  $\gamma_n(\{f \geq M_f\}) = \gamma_n(\{f \leq M_f\}) = 1/2$ , and where  $\|f\|_{\text{Lip}} = \sup_{x \neq y} |f(x) - f(y)|/\|x - y\|$ .

It has been discovered by V. Milman [M1] that (1.7) (or, equivalently, a corresponding inequality on the sphere  $S_n$ ) is at the root of the celebrated Dvoretzky's theorem. Actually, the following inequality is sufficient to prove Dvoretzky's theorem: There is a numerical constant  $K$  such that if  $f$  is a Lipschitz function on  $\mathbb{R}^n$ , then

$$\gamma_n(\{|f - \int f d\gamma_n| \geq u\}) \leq 2 \exp\left(-\frac{u^2}{K\|f\|_{\text{Lip}}^2}\right). \quad (1.8)$$

A very simple proof of this inequality (1.8) was discovered by B. Maurey and G. Pisier (cf. [P]; in that same reference is included a different proof due to Maurey using stochastic integrals which yields  $K = 2$ ).

To understand better the relationship between (1.7) and (1.8) one should note that either of these inequalities imply the fact that  $|M_f - \int f d\gamma_n| \leq K\|f\|_{\text{Lip}}$ . Here, as in the sequel,  $K$  denotes a universal constant, not necessarily the same at each occurrence. It is not our purpose here to enter the topic of local theory of Banach spaces, that was covered by Milman's paper [M2], and we turn towards the application of (1.6) to probability theory. The importance of (1.6) stems from the fact that  $\gamma$  is the prototype for all Gaussian measures. To stress the point, we now outline the “canonical” way to look at Gaussian processes, that was put forward in [D] and that turned out to be of crucial importance. Given a point  $t$  in  $\ell^2$ , the series  $\sum_{k \geq 1} t_k x_k$  converges  $\gamma$  a.e. (since  $(x_k)$  is a sequence of independent r.v.) and thereby defines an element  $X_t$  of  $L^2(\gamma)$ , of law  $N(0, \|t\|_2)$ . Any subset  $T$  of  $\ell^2$  thus defines a Gaussian process  $(X_t)_{t \in T}$ . For many purposes all Gaussian processes can be reduced to this type. We say that the process is bounded if  $\sup_{t \in T} X_t < \infty$   $\gamma$  a.e. (to avoid technicalities, we assume from now on that  $T$  is countable).

A problem of historical importance was, given a Gaussian process (that is, in our setting a subset  $T$  of  $\ell^2$ ), to understand, under the conditions that  $G$  is bounded, what are the tails of  $Y = \sup_{t \in T} |X_t|$ , i.e. the behavior of the function  $\gamma(\{Y \geq u\})$  as  $u \rightarrow \infty$ . It was proved by Landau and Shepp [L-S] and, independently by Fernique [F], that  $E(e^{\alpha Y^2}) < \infty$  for some  $\alpha > 0$ , where for simplicity, we write  $Ef$  for  $\int f d\gamma$ . Interestingly, the proof of Landau and Shepp is isoperimetric in nature. In [B1], C. Borell use (1.5) as follows. Set  $\sigma = \sup_{t \in T} \|t\|_2$ . It is then clear that  $Y(x) \leq Y(y) + \sigma u$  if  $x \in y + uB$ . Thus by (1.5)

$$\gamma(\{Y \geq M_Y + \sigma u\}) \leq \gamma_1([u, \infty)).$$

This implies that

$$\alpha > \sigma \Rightarrow E \exp \frac{Y^2}{2\alpha^2} < \infty. \quad (1.9)$$

C. Borell also used the same approach to obtain sharp integrability results for homogeneous chaos [B2, B3].

It turns out that, when more information is available on  $T$  (e.g. information about entropy numbers) results sharper than (1.10) can be obtained by specific methods. This has unfortunately lead some researchers to doubt the power of (1.6); the issue is that the usefulness of (1.6) is greatly enhanced by an appropriate use of  $A$ . This point was brought in particular to light in [T1], where the following is proved. Given a bounded process  $T \subset \ell^2$ , set

$$\tau = \inf\{u > 0; \gamma(\{\sup_{t \in T} |X_t| < u\}) > 0\}.$$

Then

$$\tau' > \tau \Rightarrow E \exp \frac{1}{2\sigma^2} (Y - \tau')^2 < \infty. \quad (1.10)$$

This result should be compared to (1.9). It can be interpreted as a tail estimate. It means that the function  $f(u) = \Phi^{-1}(\gamma(\{y \leq u\}))$  (that is concave by (1.3)) satisfies

$$0 \geq \lim_{u \rightarrow \infty} \left( f(u) - \frac{u}{\sigma} \right) \geq -\frac{\tau}{\sigma}. \quad (1.11)$$

Thus,  $f(u)$  has an asymptote  $(u/\sigma) + \ell$  with  $-\tau/\sigma \leq \ell \leq 0$ . This result is optimal in the sense that  $f$  can approach this asymptote arbitrarily slowly. We refer to [L-T2], Chapter 3, for an extension of this result to homogeneous chaos, and to [G-K] for further developments of the same idea.

While (1.10) is optimal for general processes, it can be improved when one has more information about  $T$ . In [T3] a method was introduced relying on (1.6) to improve the tail estimate (1.10) in the specific case where  $T$  is compact and there is a unique  $t \in T$  with  $\|t\| = \sigma$ . The method has been developed further in [D-M-W]. It could also be used in many other situations, e.g. to simplify the results of [B-K].

While (1.10) uses in a rather precise form the information provided by (1.6), it is often sufficient (e.g. for the proof of Dvoretzky's theorem) to have a weaker information of the type

$$\gamma(A) \geq 1/2 \Rightarrow \gamma(A + uB) \geq 1 - K \exp\left(-\frac{u^2}{K}\right) \quad (1.12)$$

without precise information on the constant  $K$ . It is this principle, rather than (1.5) that we now call the concentration of measure phenomenon (for the Gauss measure).

## 2. The Concentration of Measure Phenomenon

It seems rather unlikely that (1.6) could at all be improved, but it could come as a surprise that on the other hand (1.12) can be improved, in the sense that a similar inequality holds when the set  $A + uB$  is replaced by a smaller (and, in some cases, much smaller) set. The central result of this section is that, in the class of product measures, the natural setting for the concentration of measure phenomenon is not the Gaussian measure  $\gamma$  but rather the product measure  $v$  on  $\mathbb{R}^{\mathbb{N}}$  obtained by providing each factor with the measure  $v_1$  of density  $\frac{1}{2}e^{-|x|}$  with respect to Lebesgue measure. We set

$$B_1 = \left\{ x \in \mathbb{R}^{\mathbb{N}}; \sum |x_k| \leq 1 \right\}; B_2 = \left\{ x \in \mathbb{R}^{\mathbb{N}}; \sum x_k^2 \leq 1 \right\}.$$

**Theorem 2.1** [T6]. *There exists a universal constant  $K$  such that for all subsets  $A$  of  $\mathbb{R}^{\mathbb{N}}$ , all  $u \geq 0$ , we have*

$$v(A) = v_1((-\infty, a]) \Rightarrow v_*(A + \sqrt{u}B_2 + uB_1) \geq v_1\left(\left(-\infty, a + \frac{u}{K}\right]\right). \quad (2.1)$$

In particular

$$v(A) \geq 1/2 \Rightarrow v_*(A + \sqrt{u}B_2 + uB_1) \geq v_1\left(\left(\infty, \frac{u}{K}\right]\right) = 1 - \frac{1}{2} \exp\left(-\frac{u}{K}\right). \quad (2.2)$$

A striking difference between this inequality and (1.6) is that the set  $A$  is enlarged by the mixture  $\sqrt{u}B_2 + uB_1$  of the  $\ell^2$  and  $\ell^1$  balls, whose shape changes with the value of  $u$ . To understand the reason for the strange set  $\sqrt{u}B_2 + uB_1$ , it is instructive to derive from (2.1) the size of the tails  $v(\{X_t \geq u\})$ , where  $X_t(x) = \sum t_k x_k$  and  $t \in \ell^2$  (these can of course be obtained by a simple direct argument). The set  $A = \{X_t \leq 0\}$  satisfies  $v(A) \geq 1/2$  by symmetry. Thus by (2.1), we have  $v_*(A + \sqrt{u}B_2 + uB_1) \geq 1 - \frac{1}{2}e^{-u/K}$ . But obviously  $X_t \leq \sqrt{u}\|t\|_2 + u\|t\|_{\infty}$  on  $A + \sqrt{u}B_2 + uB_1$  (where  $\|t\|_{\infty} = \sup_{k \geq 1} |t_k|$ ). Thus we get

$$v(\{X_t \geq \sqrt{u}\|t\|_2 + u\|t\|_{\infty}\}) \leq \frac{1}{2}e^{-u/K}$$

which can be formulated as

$$v(\{X_t \geq u\}) \leq \frac{1}{2} \exp\left(-\min\left(\frac{u^2}{K\|t\|_2^2}, \frac{u}{K\|t\|_{\infty}}\right)\right)$$

(and gives the correct order for  $-\log v(\{X_t \geq u\})$ ).

Another difference between (2.1) and (1.6) is the unspecified constant  $K$  on the left side, that actually makes (2.1) closer to (1.12) than to (1.6). An interesting problem would be to find an “exact” version of (2.1). One could ask for example if there is a natural “smallest” set  $W(u)$  (whose shape would depend on  $u$  in a possibly complex way) that could be used instead of  $\sqrt{u}B_2 + uB_1$  in (2.1). The resulting inequality should give sharp estimates for  $v(\{X_t \geq u\})$ ; the variety of

competing estimates for this quantity [Hoe] might indicate the difficulty of the task.

The proof of Theorem 2.1 is made complicated by the fact that, in contrast with the Gauss measure or Lebesgue measure, the measure  $\nu$  has less symmetries (in particular is not invariant by rotations) and thus that this restricts the use of rearrangements. The method of proof is to consider a statement similar to (2.1) (the set  $\sqrt{u}B_2 + uB_1$  being replaced by a more amenable set  $C(u)$  of comparable size) and prove it by induction over  $n$ , when the set  $A$  is assumed to depend on  $n$  coordinates only. The key observation is that the proof of the induction step can be deduced from a two-dimensional statement. While the proof in (2.1) is not simple, it is beyond doubt that the important part of (2.1) is (2.2) for large values of  $u$  ( $u \geq K$ ). Fortunately, this is much simpler to prove. The idea is to prove, again by induction over the number of coordinates of which  $A$  depends, that, if one sets

$$h_A(x) = \inf\{u \geq 0; x \in A + C(u)\}$$

then  $E \exp(h_A(x)/K) \leq 1/P(A)$ , so that, by Chebyshev inequality,

$$\nu(A + C(u)) \geq 1 - \frac{1}{P(A)} \exp\left(-\frac{u}{K}\right),$$

that recovers (2.2) for  $u$  large enough.

We now explain why (2.2) is an improvement over (1.12). The argument that we will present will be referred to in the sequel as the “contraction argument”. The precise form we use was introduced by G. Pisier [P, Ch. 2] and played an essential role in the discovery of the correct formulation of Theorem 2.1. (A similar idea occurs earlier in [G-M], Section 2-1).

Consider the increasing map  $\psi$  from  $\mathbb{R}$  to  $\mathbb{R}$  that transforms  $\nu_1$  into  $\gamma_1$ . It is a simple matter to see that

$$|\psi(x) - \psi(y)| \leq K \min(|x - y|, |x - y|^{1/2}). \quad (2.3)$$

Consider the map  $\Psi$  from  $\mathbb{R}^{\mathbb{N}}$  to  $\mathbb{R}^{\mathbb{N}}$ , such that  $\Psi((x_k)_{k \geq 1}) = (\psi(x_k))_{k \geq 1}$ . Thus  $\Psi$  transforms  $\nu$  into  $\gamma$ .

Consider a Borel set  $A$  of  $\mathbb{R}^{\mathbb{N}}$  such that  $\gamma(A) \geq 1/2$ . Then

$$\begin{aligned} \gamma(\Psi(\Psi^{-1}(A) + \sqrt{u}B_2 + uB_1)) &= \nu(\Psi^{-1}(A) + \sqrt{u}B_2 + uB_1) \\ &\geq 1 - \frac{1}{2} \exp\left(-\frac{u}{K}\right). \end{aligned} \quad (2.4)$$

However, it follows from (2.3) that

$$A_u = \Psi(\Psi^{-1}(A) + \sqrt{u}B_2 + uB_1) \subset A + K\sqrt{u}B_2 \quad (2.5)$$

and thus (2.4) improves over (1.12). To illustrate the improvement of (2.4) over (1.12), consider the case where  $A = \{x; \forall k \leq n, |x_k| \leq a_n\}$ , where  $a_n$  is chosen so that  $\gamma(A) = 1/2$  (and hence is of order  $\sqrt{\log n}$ ). Then, for  $u \ll \log n$ , the set  $A_u$  is easily seen to be contained in

$$A + K \left( \sqrt{\frac{u}{\log n}} B_2 + \frac{u B_1}{\sqrt{\log n}} \right) \subset A + \left( \frac{K u}{\log n} \right)^{1/2} \sqrt{u} B_2,$$

where  $u/\log n \ll 1$ .

One intriguing aspect of Theorem 2.1, when used as an improvement over (1.12), is that it breaks the rotational invariance of the Gauss measure. Indeed, it not only tells us that  $\gamma_n(A_u) \geq 1 - \exp(-u/K)$  (where  $A_u$  is defined in (2.5)) but also that  $\gamma_n((RA)_u) \geq 1 - \exp(-u/K)$  for any rotation  $R$  of  $\mathbb{R}^n$ .

A natural question is whether (2.2) is the correct formulation of the concentration of measure phenomenon. This seems to be the case, at least in the setting of product measures. Indeed, consider a probability measure  $\theta_1$ , on  $\mathbb{R}$ , and its product  $\theta$  on  $\mathbb{R}^N$ . Suppose that the following holds (that is much weaker than (1.2)). There exists  $K > 0$ , such that

$$\theta(A) \geq 1/2 \Rightarrow \theta(A + KB_\infty) \geq 3/4$$

where  $B_\infty = \{x \in \mathbb{R}^N, \forall k \geq 1, |x_k| \leq 1\}$ . Then the tails  $f(u) = \theta(\{|x| \geq u\})$  must decay exponentially [T6]. Note that, if these tails decay exponentially in a smooth enough way,  $\theta_1$  is the image of  $v_1$  by a contraction, and Pisier's contraction argument presented before shows that (2.1) will also hold for  $\theta$ .

Consider now  $1 \leq \alpha < \infty$  and the measure  $v^\alpha$  on  $\mathbb{R}^N$ , obtained as the product measure when each factor is endowed with the probability measure  $a_\alpha e^{-|x|^\alpha} dx$  (where  $a_\alpha$  is a normalizing constant). The contraction argument presented above shows that

$$v^\alpha(A) \geq 1/2 \Rightarrow v_*^\alpha(A + U_\alpha(u)) \geq 1 - \exp\left(-\frac{u}{K(\alpha)}\right). \quad (2.6)_\alpha$$

where  $U_\alpha(u) = u^{1/2} B_2 + u^{1/\alpha} B_\alpha$  for  $\alpha \leq 2$ ,  $U_\alpha(u) = u^{1/2} B_2 \cap u^{1/\alpha} B_\alpha$  for  $\alpha \geq 2$ , and  $B_\alpha = \{x \in \mathbb{R}^N; \sum |x_k|^\alpha \leq 1\}$ . For  $\alpha > \beta$ ,  $(2.6)_\alpha$  is a consequence of  $(2.6)_\beta$  (by the contraction argument).

As in the Gaussian case, to each point  $t \in \ell^2$  one can associate the random variable  $X_t = \sum_{k \geq 1} t_k x_k$  on  $(\mathbb{R}^N, v^\alpha)$ ; and each subset  $T$  of  $\ell^2$  thus defines a stochastic process. The main motivation for proving  $(2.6)_\alpha$  was the discovery [T7, T8] of a new approach to the problem of finding lower bounds for  $E \sup_{t \in T} X_t$  that makes  $(2.6)_\alpha$  an essential step. This new approach eliminates the use of Slepian's lemma [S], which is a specific property of Gaussian processes. It replaces it by the use of  $(2.6)_\alpha$ , combined with a Sudakov-type minoration [Su]. It enables to describe  $E \sup_{t \in T} X_t$  in terms of the geometry of  $T$ , thereby extending the results of [T2] for the Gaussian case  $\alpha = 2$ . But due to limitations of space we cannot discuss this point further.

### 3. Concentration of Measure for Bernoulli Random Variables

Pisier used his contraction argument mentioned above to conclude from (1.2) that if  $\lambda_n$  denotes the product measure on  $\mathbb{R}^n$  when  $\mathbb{R}$  is equipped with the uniform measure  $\lambda_1$  on  $[-1, 1]$ , then

$$\lambda_n(A) = \Phi(a) \Rightarrow \lambda_n(A_u) \geq \Phi\left(a + \frac{u}{K}\right). \quad (3.1)$$

Closely related to  $\lambda_n$ , but of somewhat greater importance in Probability (since it corresponds to random signs) is the probability  $\mu_n$  on  $\{-1, 1\}^n$  that gives mass  $2^{-n}$  to each point. The problem arises whether a concentration of measure principle as strong as (3.1) holds for  $\mu_n$ . This is not the case (as follows from the example given after (3.3)). The appropriate formulation for a substitute to (3.1) requires to think to  $\{-1, 1\}^n$  as a subset of  $\mathbb{R}^n$ . For a non-empty subset  $A$  of  $\{-1, 1\}^n$ , we set  $\varphi_A(x) = \inf\{\|x - y\|_2; y \in \text{conv } A\}$ , where  $\text{conv } A$  denotes the convex hull of  $A$  in  $\mathbb{R}^n$ .

**Theorem 3.1** [T4].  $E \exp(\varphi_A^2/8) \leq 1/\mu_n(A)$ .

Using Chebyshev inequality gives

$$\text{For } u \geq 0, \mu_n(\{\varphi_A \geq u\}) \leq \frac{1}{\mu_n(A)} e^{-u^2/8}. \quad (3.2)$$

We first explore the consequences of this result. Consider a *convex* function for  $\mathbb{R}^n$ . Then one can deduce from (3.2) that if  $M_f$  is median of  $f$  (for  $\mu_n$ ), we have

$$\mu_n(\{|f - M_f| > u\}) \leq 4 \exp\left(-\frac{u^2}{8\|f\|_{\text{Lip}}^2}\right). \quad (3.3)$$

This inequality should be compared to (1.7). A major difference with (1.7) is however that this result is really specific to *convex* functions. To see it, consider  $n$  even, and let  $A = \{x \in \{-1, 1\}^n; \sum_{i \leq n} x_i \leq 0\}$ , so that  $\mu_n(A) \geq 1/2$ . Define  $f(x) = \inf\{\|x - y\|_2 : y \in A\}$ , so that  $\|f\|_{\text{Lip}} = 1$ , and  $M_f = 0$ . It is easy to see that for  $y \in \{-1, 1\}^n$ , we have  $f(y) = \sqrt{2}((\sum_{i \leq n} y_i)^+)^{1/2}$ . But the central limit theorem shows that  $\mu_n(\{f \geq cn^{1/4}\}) \geq 1/4$  for some constant  $c$  independent of  $n$ . (Note then that  $\mu_n(A_{cn^{1/4}}) \leq 3/4$  and that (3.1) fails.)

Consider now a Banach space  $E$  and vectors  $(x_k)_{k \leq n}$  in  $E$ .

Set

$$\sigma = \sup \left\{ \sum_{k \leq n} x^*(x_k)^2; x^* \in E^*, \|x^*\| \leq 1 \right\}.$$

The function on  $\mathbb{R}^n$  given by  $f(y) = \|\sum_{k \leq n} y_k x_k\|_E$  is convex and satisfies  $\|f\|_{\text{Lip}} = \sigma$ . Consider a sequence  $(\varepsilon_k)_{k \leq n}$  of (symmetric) Bernoulli random variables; that is, the sequence is independent identically distributed and  $P(\varepsilon_i = 1) = P(\varepsilon_i = -1) = 1/2$ . Then (3.3) implies

$$P \left( \left| \left\| \sum_{k \leq n} \varepsilon_k x_k \right\| - M \right| \geq u \right) \leq 4e^{-u^2/8\sigma^2} \quad (3.4)$$

where  $M$  is a median of  $\|\sum_{k \leq n} \varepsilon_k x_k\|$ . From (3.4) and elementary computations follows that

$$\left\| \sum_{k \leq n} \varepsilon_k x_k \right\|_p \leq \left\| \sum_{k \leq n} \varepsilon_k x_k \right\|_1 + K\sigma p^{1/2},$$

a precise form of the so called Khintchine-Kahane inequalities.

It is not known whether the exponent  $1/8$  in (3.2) can be improved; the best possible exponent would be  $1/2$ . Another problem of interest would be the determination of  $\min\{\mu_n(\{\varphi_A \geq u\}); \mu_n(A) = u\}$ . It is likely that the sets which achieve this minimum depend on  $A, u$ ; thus the problem might be difficult.

It is of interest to compare Theorem 3.1 with the classical results concerning Hamming distance. The Hamming distance  $d(s, t)$  of two points  $s, t$  of a product of sets is the number of coordinates where  $s, t$  differ. For a subset  $A$  of  $\{-1, 1\}^n$ , we set  $d_A(x) = \inf\{d(x, y), y \in A\}$ . It follows from an isoperimetric inequality of Harper [Ha] that for  $\mu_n(A) \geq 1/2$ , we have

$$\mu_n(\{d_A \geq u\sqrt{n}\}) \leq \exp(-2u^2). \quad (3.5)$$

On the other hand, it is simple to see that  $2d_A \leq \sqrt{n}\varphi_A$ . Thus  $\{d_A \geq u\sqrt{n}\} \subset \{\varphi_A \geq 2u\}$ . Now (3.2) provides the estimate

$$\mu_n(\{\varphi_A \geq 2u\}) \leq 2 \exp(-u^2/2).$$

Compared with (3.5), this provides a weaker bound (but of the same essential strength) for a larger set. The most important difference is however that (3.2), in contrast with (3.5), is independent of the dimension.

We now present an “abstract” extension of Theorem 3.1. Consider a sequence  $(\Omega_k, \mu_k)_{k \leq n}$  of probability spaces and denote by  $P$  the product measure on  $\Omega = \prod_{k \leq n} \Omega_k$ . Consider a subset  $A$  of  $\Omega$ . For  $x \in \Omega$ , consider the set

$$U_A(x) = \{t \in \{0, 1\}^n; \exists y \in A, t_k = 0 \Rightarrow x_k = y_k\}.$$

We consider  $U_A(x)$  as a subset of  $\mathbb{R}^n$ ; we denote by  $V_A(x)$  the convex hull of  $U_A(x)$ .

For  $\alpha \geq 1$ ,  $0 \leq u \leq 1$ , we consider the function

$$\xi(\alpha, u) = \alpha(1-u) \log(1-u) - (\alpha+1-\alpha u) \log\left(1 - \frac{\alpha u}{1+\alpha}\right).$$

Elementary calculus show that this function increases in  $\alpha, u$ , and is convex in  $u$ .

We set

$$\psi_{\alpha, A}(x) = \inf \left\{ \sum_{i \leq n} \xi(\alpha, y_i); y = (y_i)_{i \leq n} \in V_A(x) \right\}.$$

**Theorem 3.2** [T9].  $E \exp \psi_{\alpha, A} \leq (1/P_*(A))^\alpha$ .

Calculus shows that  $\xi(1, u) \geq u^2/4$ ; thus Theorem 3.2 implies Theorem 3.1, but only with the worse exponent  $1/16$  instead of  $1/8$ . An essential improvement of Theorem 3.2 over Theorem 3.1 is that for  $\alpha$  large and  $u$  close to 1,  $\xi(\alpha, u)$  is of order  $\xi(\alpha, 1) = \log(1 + \alpha)$ . The following bound seems to be of particular interest. For  $t \geq 1$ ,

$$P_x(A) \geq 1/2 \Rightarrow P(\{\psi_{t,A} \geq t\}) \leq (2/e)^t.$$

It has been observed in [J-S] (using the method of [T4]) that if  $0 < \eta < 1$ , and if  $\mu_n$  denotes now the measure  $((1 - \eta)\delta_0 + \eta\delta_1)^n$  on  $\{0, 1\}^n$ , then for a set  $A \subset \{0, 1\}^n$ , we have  $E \exp\{\varphi_A^2/4\} \leq 1/\mu_n(A)$  (this also follows from Theorem 3.2). An interesting fact in that direction is that the tails of  $\varphi_A$  do not improve when  $\eta$  is small. This is somewhat unexpected. To see it, consider the case where  $A = \{x \in \{0, 1\}^n; \sum_{k \leq n} x_k \leq \eta n\}$ , so that  $\mu_n(A)$  is of order  $1/2$  by the law of large numbers. On the other hand, it is simple to see that (for  $\eta n$  integer)  $\varphi_A(y) \leq u$  if and only if  $\sum y_k = p$  where  $\sqrt{p}(1 - \eta n/p) \leq u$ , so that  $p \leq \eta n + u\sqrt{p}$ . For  $u \leq (\eta n)^{1/2}$ , this implies  $p \leq 2\eta n$ , so that  $p \leq \eta n + u\sqrt{2\eta n}$ . Thus for  $p > \eta n + u\sqrt{2\eta n}$  we have  $\varphi_A(y) > u$ . It follows from the central limit theorem that if  $0 < \eta \leq 1/2$ , then for  $n$  large enough, we have  $\mu_n(\{\varphi_A(y) > u\}) \geq \exp(-cu^2)$  for some  $c$  independent of  $n, \eta$ .

#### 4. An Isoperimetric Inequality for Product Measure

An important concentration of measure phenomenon for product measures is as follows. Consider a sequence  $(\Omega_k, \mu_k)_{k \leq n}$  of probability spaces. Denote by  $P$  the product measure on  $\Omega = \prod_{k \leq n} \Omega_k$ . Then

$$P(A) \geq 1/2 \Rightarrow P(\{d_A \geq u\}) \leq 2 \exp(-u^2/Kn). \quad (4.1)$$

where the Hamming distance  $d_A$  has been introduced in Section 3. This is an extension of (3.5) (with worse constants). It is easy to prove using the martingale approach introduced by B. Maurey and developed by G. Schechtman (see [M-S]). It also follows from Theorem 3.2 the way (3.5) follows from Theorem 3.1. (This approach gives a constant  $K = 4$  in the exponent.)

For a set  $A \subset \Omega$ , and  $k, q \geq 0$ , consider

$$H(A, q, k) = \left\{ y \in \prod_{k \leq n} \Omega_k; \exists x^1, \dots, x^q \in A; \text{ card } \{i; \forall \ell \leq q, y_i \neq x_i^\ell\} \leq k \right\}.$$

For  $q = 1$ , this is exactly the set  $\{d_A \leq k\}$ . The set  $H(A, q, k)$  can be thought of as an ‘‘enlargement’’ of  $A$ , although it does not seem possible to define it as a neighborhood of  $A$  for a distance.

**Theorem 4.1** [T5]. *For some universal constant  $K$ , and all  $k, q \geq 1$ , we have*

$$P(A) \geq 1/2 \Rightarrow P_*(H(A, q, k)) \geq 1 - \left( \frac{K}{k} + \frac{K}{q \log q} \right)^k. \quad (4.2)$$

As stated, this theorem gives information only when  $k, q$  are large. However it is also possible to show that if  $q \geq 2, k \geq k_0$ , then

$$P(A) \geq 1/2 \Rightarrow P_*(H(A, q, k)) \geq 1 - \eta^k$$

where  $\eta < 1$  is universal. In contrast with the case  $q = 1$  (4.1), the estimate (4.2) is independent of the number of coordinates (and thus can be extended to the case of an infinity of factors.)

To gain some intuition about (4.2), it is useful to consider the case where  $\Omega_i = \{0, 1\}, \mu_i(\{0\}) = 1 - 1/n, \mu_i(\{1\}) = 1/n$ , and

$$A = \left\{ (x_i); \sum_{i \leq n} x_i \leq 1 \right\}.$$

In that case,  $P(A) \geq 1/2$  and

$$H(A, q, k) = \left\{ (x_i); \sum_{i \leq n} x_i \leq q + k \right\}.$$

For  $k$  of order  $q \log q$ , simple estimates show that

$$P(H(A, q, k)) \leq 1 - \left( \frac{1}{Kq \log q} \right)^k,$$

which should be compared to (4.2).

Theorem 4.1 has strong implications about the behavior of sums of independent random variables valued in a Banach space. Consider such variables  $X_1, \dots, X_n$  valued into a separable Banach space  $F$ . We now outline a method to obtain bounds on the tails of  $\|\sum_{i \leq n} X_i\|$ . (These bounds can now also be derived from Theorem 3.2, which has a considerably simpler proof than Theorem 4.1. Tail estimates are in particular at the root of classical theorems like strong laws of large numbers and laws of the iterated logarithm.) While this method might look complicated at first glance, it seems to capture the size of these tails in essentially all the situations; see e.g. [T5, L-T1]. Without essential loss of generality, one can assume that the variables are symmetric, i.e.  $X_i$  has the same distribution as  $-X_i$ . Consider a sequence  $(\varepsilon_i)_{i \leq n}$  of Bernoulli random variables, that can be assumed to be independent of  $(X_i)_{i \leq n}$ . Thus  $\sum_{i \leq n} \varepsilon_i X_i$  has the same distribution as  $\sum_{i \leq n} X_i$ . We then write

$$\begin{aligned} \left\| \sum_{i \leq n} \varepsilon_i X_i \right\| &= E_\varepsilon \left\| \sum_{i \leq n} \varepsilon_i X_i \right\| + \left( \left\| \sum_{i \leq n} \varepsilon_i X_i \right\| - E_\varepsilon \left\| \sum_{i \leq n} \varepsilon_i X_i \right\| \right) \\ &:= (\text{I}) + (\text{II}) \end{aligned} \tag{4.3}$$

where  $E_\varepsilon$  is the conditional expectation given  $(X_i)_{i \leq n}$ . Denote by  $\mu_i$  the law of  $X_i$  on  $F$ ; Consider a set  $A \subset F^n$ , and suppose that

$$A \subset \left\{ (x_1, \dots, x_n) \in F^n, E_\varepsilon \left\| \sum_{i \leq n} \varepsilon_i x_i \right\| \leq M \right\}. \tag{4.4}$$

Then it is easy to see that

$$(x_1, \dots, x_n) \in H(A, q, k) \Rightarrow \left\| \sum_{i \leq n} \varepsilon_i x_i \right\| \leq qM + \sum_{i \leq k} \|x_i\|^*$$

where  $\|x_i\|^*$  is the  $i$ -th largest term of the sequence  $(\|x_i\|)_{i \leq n}$ . Suppose now that  $P((X_1, \dots, X_n) \in A) \geq 1/2$  (e.g. if  $M = 2E\|\sum_{i \leq n} \varepsilon_i X_i\|$  and there is equality in (4.4)). It then follows from (4.2) that, if  $k \geq q$

$$P((I) \geq qM + \sum_{i \leq k} \|X_i\|^*) \leq \left(\frac{K}{q}\right)^k. \quad (4.5)$$

On the other hand, if we set

$$\sigma_X^2 = \sup \left\{ \sum_{i \leq n} x^*(X_i)^2; x^* \in E^*, \|x^*\| \leq 1 \right\},$$

it follows from (3.4) that, conditionally on  $X_1, \dots, X_n$

$$P((II) \geq K(1+u)\sigma_X) \leq 4e^{-u^2}. \quad (4.6)$$

To make (4.5), (4.6) usable, it remains to control  $\sum_{i \leq n} \|X_i\|^*$  (which is a problem about real-valued random variables) and  $\sigma_X$ . Several methods have been developed for that purpose; adjusting the various parameters involved has allowed to get bounds of the right order in all the problems studied to date; cf. [L-T2], Chapters 6 to 8.

## References

- [A] R. J. Adler: An introduction to continuity, extrema, and related topics for general Gaussian processes. Forthcoming book
- [B-K] S. Berman, N. Kôno: The maximum of a Gaussian process with non constant variance; a sharp bound for the distribution tail. Ann. Probab. **17**, 632–650 (1989)
- [B1] C. Borell: The Brunn-Minkowski inequality in Gauss space. Invent. math. **30**, 207–216 (1975)
- [B2] C. Borell: Tail probabilities in Gauss space. Vector Spaces Measures and Application, Dublin 1977. (Lecture Notes in Mathematics, vol. 644.) Springer, Berlin Heidelberg New York 1978, pp. 71–82
- [B3] C. Borell: On polynomial chaos and integrability. Prob. Math. Statist. **3**, 191–203 (1984)
- [D] R. M. Dudley: The sizes of compact subsets of Hilbert space and continuity of Gaussian processes. J. Funct. Anal. **1**, 290–330 (1987)
- [D-M-W] V. Dobric, M. B. Marcus, M. Weber: The distribution of the large values of the supremum of a Gaussian process
- [E1] A. Ehrhard: Symétrisation dans l'espace de Gauss. Math. Scand. **53**, 281–301 (1983)
- [E2] A. Ehrhard: Inégalités isoperimétriques et intégrales de Dirichlet Gaussiennes. Ann. Sci. Ec. Norm. Sup. **17**, 317–332 (1984)

- [E3] A. Ehrhard: Éléments extrémaux pour les inégalités de Brunn-Minkowski gaussiennes. *Ann. Inst. H. Poincaré* **22**, 149–168 (1986)
- [F] X. Fernique: Intégrabilité des vecteurs gaussiens. *C. R. Acad. Sci. Paris* **270**, 1698–1699 (1970)
- [F-L-M] T. Figiel, J. Lindenstrauss, V. Milman: The dimensions of almost spherical sections of convex bodies. *Acta Math.* **139**, 52–94 (1977)
- [G-K] V. Goodman, J. Kuelbs: Rates of clustering of some self similar Gaussian processes. *Manuscript*, 1990
- [G-M] M. Gromov, V. D. Milman: A topological application of the isoperimetric inequality. *Amer. J. Math.* **105**, 843–854 (1983)
- [Ha] L. H. Harper: Optimal numbering and isoperimetric problems on graphs. *J. Comb. Theory* **1**, 385–393 (1966)
- [Hoe] W. Hoeffding: Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58**, 13–30 (1963)
- [Hof] J. Hoffmann-Jørgensen: Sums of independent Banach space valued random variables. *Studia Math.* **52**, 159–186 (1974)
- [J-S] W. B. Johnson, G. Schechtman: Remarks on Talagrand's deviation inequality for Rademacher functions. *Texas Functional Analysis Seminar 1988–89*. The University of Texas, 1989
- [K] J. P. Kahane: Some random series of functions. Cambridge University Press, 1985
- [K-S] S. Kwapień, J. Szulga: Hypercontraction methods for comparison of moments of random series in normed spaces. *Ann. Probab.* (to appear)
- [L-S] H. Landau, L. A. Shepp: On the supremum of a Gaussian process. *Sankhyā* **32**, 369–378 (1970)
- [L] M. Ledoux: Isoperimétrie et inégalités de Sobolev logarithmiques gaussiennes. *C. R. Acad. Sci. Paris* **306**, 79–82 (1988)
- [Lé] P. Lévy: Problèmes concrets d'analyse fonctionnelle. Gauthier Villars, 1951
- [L-T1] M. Ledoux, M. Talagrand: Some application of isoperimetric methods to strong limit theorems for sums of independent random variables, *Ann. Probab.* **18**, 754–789 (1990)
- [L-T2] M. Ledoux, M. Talagrand: Probability in Banach spaces. (Ergebnisse der Mathematik, Bd. 23). Springer, Berlin Heidelberg New York 1991
- [McK] H. P. McKean: Geometry of differential space. *Ann. Probab.* **1**, 197–206 (1973)
- [M1] V. D. Milman: New proof of the theorem of Dvoretzky on sections of convex bodies. *Funct. Anal. Appl.* **5**, 28–37 (1971)
- [M2] V. D. Milman: The Concentration Phenomenon and Linear Structure of Finite Dimensional Normed Spaces. Proceedings of the International Congress of Mathematicians, Berkeley 1986, pp. 961–975
- [M3] V. D. Milman: The heritage of P. Lévy in functional analysis. *Asterisque* **157–158**, 273–302 (1988)
- [M-S] V. D. Milman, G. Schechtman: Asymptotic theory of finite dimensional normed spaces. (Lecture Notes in Mathematics, vol. 1200.) Springer, Berlin Heidelberg New York 1986
- [O] R. Osserman: The isoperimetric inequality. *Bull. Amer. Math. Soc.* **84**, 1182–1238 (1978)
- [P] G. Pisier: Probabilistic methods in the geometry of Banach spaces. *Probability and Analysis, Varenna (Italy) 1985*. (Lecture Notes in Mathematics, vol. 1206.) Springer, Berlin Heidelberg New York 1986, pp. 167–241
- [S] D. Slepian: The one sided barrier problem for Gaussian noise. *Bell System Tech. J.* **41**, 463–501 (1962)

- 
- [Su] V. N. Sudakov: Gaussian measures, Cauchy measures and  $\varepsilon$ -entropy. Sov. Math. Dokl. **10**, 310–315 (1969)
  - [S-T] V. N. Sudakov, B. S. Tsirelson: Extremal properties of half spaces for spherically invariant measures. J. Sov. Math. **9**, 9–18 (1978) [translated from Zap. Nauch. Sem. L.O.M.I. **41**, 14–24 (1974)]
  - [T1] M. Talagrand: Sur l'intégrabilité des vecteurs gaussiens. Z. Wahrscheinlichkeitstheorie verw. Gebiete **68**, 1–8 (1984)
  - [T2] M. Talagrand: Regularity of Gaussian processes. Acta Math. **159**, 99–149 (1987)
  - [T3] M. Talagrand: Small tails for the supremum of a Gaussian process. Ann. Inst. H. Poincaré **24**, 307–315 (1988)
  - [T4] M. Talagrand: An isoperimetric theorem on the cube and the Kintchine-Kahane inequalities. Proc. Amer. Math. Soc. **104**, 905–909 (1988)
  - [T5] M. Talagrand: Isoperimetry and integrability of the sum of independent Banach-space valued random variables. Ann. Probab. **17**, 1546–1570 (1989)
  - [T6] M. Talagrand: A new isoperimetric inequality and the concentration of measure phenomenon, to appear in Geometric aspects of functional analysis (Israel seminar). (Lecture Notes in Mathematics, vol. 1469.) Springer, Berlin Heidelberg New York 1991
  - [T7] M. Talagrand: Supremum of canonical processes. Manuscript, 1990
  - [T8] M. Talagrand: Regularity of infinitely divisible processes. Manuscript, 1990
  - [T9] M. Talagrand: A new isoperimetric inequality for product measure and the tails of sums of independent random variables. Manuscript, 1990

# Random Walks and Diffusions on Fractals

Martin T. Barlow

Statistical Laboratory, University of Cambridge, 16 Mill Lane  
Cambridge CB2 1SB, England

## 1. Introduction

The field this paper will survey could be called *analysis on fractals*: more specifically it is the study (by analytic or probabilistic means) of the fundamental second order PDEs on fractal spaces. The original motivation came from mathematical physicists working on the properties of disordered media, and interested in questions such as heat conduction, vibration modes etc. There is experimental evidence that fractals can provide a good model for certain kinds of disordered media; and it is hoped that the study of PDEs on fractal spaces would give at least some information about PDEs in disordered media. See [AO, RT] for early work by mathematical physicists, and [HBA] for a survey of the now very extensive physics literature.

Let  $F \subseteq \mathbb{R}^d$  be a connected self-similar fractal with Hausdorff dimension  $d_f = d_f(F)$ , and let  $\mu_F$  denote Hausdorff  $x^{d_f}$ -measure on  $F$ . (See Hutchinson [H] for a construction of such sets via families of linear maps.) The heat equation on  $F$  should take the form

$$\Delta_F u = \frac{\partial u}{\partial t}, \quad u(x, 0) = u_0(x), \quad x \in F, \quad (1.1)$$

where  $u : F \times \mathbb{R}_+ \rightarrow \mathbb{R}$ ,  $u_0 \in C(F)$ , and  $\Delta_F$  is a ‘Laplacian’ operator (i.e. self-adjoint with respect to  $\mu$ , local, non-positive, satisfying  $\Delta_F 1 = 0$ ) acting on a subspace  $\mathcal{D}(\Delta_F) \subset C(F)$ . The following problems arise immediately:

- (a) *Existence.* The construction of a suitable  $\Delta_F$  which is  $F$ -isotropic, that is, locally invariant with respect to local isometries of  $F$ .
- (b) *Uniqueness.* Is  $\Delta_F$  characterised (up to a scale factor) by the property of being  $F$ -isotropic?
- (c) *Properties.* The form of the functions in  $\mathcal{D}(\Delta_F)$ , and the properties of solutions to the Laplace, heat etc. equations associated with  $F$ .

Although the questions stated above are purely analytic, a natural approach (and historically the first one) is to approach them probabilistically, and seek to construct an  $F$ -isotropic diffusion process  $X_t, t \geq 0$  on  $F$ . Then the infinitesimal generator of  $X$  is a natural candidate for the Laplacian  $\Delta_F$ , and the transition density of  $X$  is a solution to the heat equation. Of course the process  $X$  is an interesting object, worthy of study in its own right.

In this survey I will concentrate on the probabilistic approach, on the construction and behaviour of diffusion processes on some specific classes of fractals. The subject is still quite young, and though some general patterns and unifying approaches are beginning to emerge, it still consists to a considerable extent of a collection of specific examples.

I will also concentrate on rigorous results in the mathematical literature. However, many of the key concepts, and in particular the indices  $d_f$ ,  $d_w$  and  $d_s$  were introduced earlier in physics papers.

## 2. The Sierpinski Gasket

This is the simplest non-trivial connected fractal, and is the natural starting point for any investigation of fractal spaces. Hutchinson's theory ([H, F]) provides a convenient description. Let  $a_1 = (0, 0)$ ,  $a_2 = (1, 0)$ ,  $a_3 = (\frac{1}{2}, \frac{1}{2}\sqrt{3})$ , let  $G_0 = \{a_1, a_2, a_3\}$ ,  $A_0$  be the closed convex hull of  $G_0$ , and let  $\phi_i$ ,  $1 \leq i \leq 3$ , be the linear contractions defined by  $\phi_i(x) = \frac{1}{2}(x + a_i)$ . For any set  $B \subseteq \mathbb{R}^2$  set

$$\Phi(B) = \bigcup_i \phi_i(B), \quad (2.1)$$

and let  $\Phi^n$  be the  $n$ -fold convolution of  $\Phi$ . The Sierpinski gasket  $G$  may be defined by

$$G = \bigcap_{n=0}^{\infty} \Phi^n(A_0) = cl \left( \bigcup_{n=0}^{\infty} \Phi^n(G_0) \right). \quad (2.2)$$

Note that  $\Phi^n(A_0)$  consists of  $3^n$  triangles each of side  $2^{-n}$ . Following the terminology of Lindström [L] we call any set of the form  $B \cap G$  (respectively  $B \cap \Phi^n(G_0)$ ) an  $n$ -complex (respectively  $n$ -cell), where  $B$  is any triangle of side  $2^{-n}$  in  $\Phi^n(A_0)$ . Write  $G_n = \Phi^n(G_0)$ .

By [H] the Hausdorff dimension of  $G$  is given by  $\dim G = \log 3 / \log 2 = d_f(G) = d_f$ . Let  $\mu_G$  be the multiple of Hausdorff- $x^{d_f}$  measure on  $G$  which assigns mass 1 to  $G$ . The key property of  $G$  used below is that it is finitely ramified:

- (FR) if  $H_1$  and  $H_2$  are two adjacent  $n$ -complexes, with common point  $x$ , then any continuous path in  $H_1 \cup H_2$  from  $H_1$  to  $H_2$  passes through  $x$ .

The natural approach to the construction of a diffusion  $X$  on  $G$  is to approximate  $X$  by a sequence of random walks on the discrete sets  $G_m$ . Call any pair of points  $x, y \in G_m$  neighbours if  $x$  and  $y$  belong to the same  $m$ -cell, and write  $N_m(x)$  for the set of neighbours of  $x$ : note that  $\#N_m(x) = 4$  for  $x \in G_m \setminus G_0$ . Let  $Y^m(r)$ ,  $r \geq 0$  be a simple symmetric nearest-neighbour random walk on  $G_m$ . Initially we will take  $Y^m(0) = 0$  for each  $m$ .

Define the sequences of stopping times

$$S_m = \min\{r \geq 0 : Y^m(r) \in G_0 \setminus \{0\}\},$$

$$T_{k+1}^m = \min\{r \geq T_k^m : Y^m(r) \in G_{m-1} \setminus \{Y^m(T_k^m)\}\}, \quad T_0^m = 0.$$

Set

$$\tilde{Y}_k^{m-1} = Y^m(T_k^m), \quad k \geq 0;$$

we call  $\tilde{Y}^{m-1}$  the *decimated random walk on  $G_{m-1}$* . The following lemma, which is crucial in the analysis of the walks  $Y^m$ , is any easy consequence of the geometry of  $G$ , and in particular of (FR).

**Lemma 1** ('Decimation invariance' [G, Ku1, BP]).

- (a) *The random walks  $Y^{m-1}$  and  $\tilde{Y}^{m-1}$  are equal in law.*
- (b) *The r.v.  $(T_{k+1}^m - T_k^m, k \geq 0)$  are i.i.d. and equal in law to  $S_1$ .*

A straightforward calculation gives  $ES_1 = 5$ , and using Lemma 1(b) it follows that  $ES_m = 5^m$ . This suggests that, to obtain a weak limit, one should consider the processes

$$X_t^{(m)} = Y^m([5^m t]) \quad t \geq 0. \quad (2.3)$$

**Theorem 2** [G, Ku1, BP]. *The processes  $S^{(m)}$ ,  $m \geq 0$ , converge weakly to a non-trivial continuous  $G$ -valued process  $X$ .*

The tightness of the sequence  $X^{(m)}$  follows from the estimates on the crossing times  $S_m$ , and similar estimates on the crossing times of smaller  $n$ -complexes. The decimation invariance of the random walks  $Y^m$  provides the necessary connection between  $Y^m$  and  $Y^{m-1}$  to prove convergence of the entire sequence  $X^{(m)}$ , rather than just a subsequence.

A similar argument can be used to define  $X$  with  $X_0 = x \in \bigcup G_n$ . Let  $X_t(w) = w(t)$  be defined on the canonical space  $\Omega = C(\mathbb{R}_+, G)$ , and let  $P^x$  be the law of  $X$  started at  $x$ . Then the law  $P^x$  for any  $x \in G$  may be defined by considering the limits  $P^{x_n}$ , as  $x_n \rightarrow x$ . Set

$$P_t f(x) = E^x f(X_t), \quad f \in C(G) \quad x \in G. \quad (2.4)$$

**Theorem 3** [Ku1, BP]. (a) *The process  $(P^x, X_t)$  is a continuous Strong Markov process on  $G$ .*

- (b)  *$P_t$  is a Feller semigroup.*
- (c)  *$P_t$  is  $\mu$ -symmetric, that is*

$$\int g(x) P_t f(x) \mu(dx) = \int P_t g(x) f(x) \mu(dx) \quad \text{for all } f, g \in C(G).$$

The proofs of (a) and (b) use no really new ideas, but do involve a number of somewhat messy technicalities. It would be nice to see a really clean construction. (c) is an immediate consequence of the symmetry of the random walks  $Y^m$ .

Let  $(\mathcal{L}, \mathcal{D}(\mathcal{L}))$  be the infinitesimal generator of the semigroup  $P_t$ : then  $\mathcal{L}$  is a Laplacian operator on  $G$ , so taking  $\Delta_G = \mathcal{L}$  gives a solution to the existence problem (2.1)(a). Write  $\mathcal{E}(f, f)$  for the Dirichlet form associated with the semigroup  $P_t$  – see [Fu].

We now consider the second problem, that is, the uniqueness of  $X$  and  $\mathcal{L}$ . A *local isometry* of  $G$  is a triple  $(H, J, \phi)$ , where  $H, J$  are subsets of  $G$ , and  $\phi : H \rightarrow J$  is an isometry in the intrinsic metric  $d_G$ . ( $d_G(x, y)$  is the Euclidean length of the shortest path between  $x$  and  $y$  in  $G$ ; it is easily seen that the intrinsic and Euclidean metrics are equivalent.) Let  $\Pi_1$  be the set of local isometries of the form  $(H, J, \phi)$  where  $H, J$  are  $n$ -complexes, and  $\phi$  is the restriction to  $G$  of a linear isometry in  $\mathbb{R}^2$ . Let  $\Pi_2$  be the set of isometries of the form  $(H_1 \cup H_2, H_1 \cup H_2, \phi)$ ,

where  $H_1, H_2$  are a pair of adjacent  $n$ -complexes, with common point  $x$  say, and  $\varphi$  fixes one of the  $H_i$  and reflects the other in the axis of symmetry which passes through  $x$ .

It is easy to see from the corresponding property of the  $Y^m$  that  $X$  is locally invariant under local isometries. By analogy with the case of  $\mathbb{R}^d$  (where standard Brownian motion is the unique diffusion invariant with respect to translations and rotations) it is therefore reasonable to call  $X$  *Brownian motion* on  $F$ .

**Theorem 4** [BP, Theorem 8.1]. *Let  $(Y(t), \tilde{P}^x)$  be a diffusion process on  $G$ , and for  $H \subset G$  write  $T_H = \inf\{t \geq 0 : Y_t \notin H\}$ . Suppose that, whenever  $(H, J, \phi) \in \Pi_1 \cup \Pi_2$  is a local isometry of  $G$ , and  $x \in H$ , the  $\tilde{P}^x$  law of  $\phi(Y(T_J \wedge \cdot))$  and the  $\tilde{P}^{\phi(x)}$  law of  $Y(T_J \wedge \cdot)$  are equal. Then there exists a constant  $a \in [0, \infty)$  such that the  $\tilde{P}^x$  law of  $(Y(t), t \geq 0)$  equals the  $P^x$  law of  $(X_{at}, t \geq 0)$ .*

*Remark.* The class of isometries required to ensure uniqueness is slightly larger than one would initially expect. The isometries in the class  $\Pi_2$  are necessary to exclude the process (similar to some discussed in [G]) which moves along the edge of the largest  $m$ -complex it lies in.

We now turn to the properties of the process  $X$ . Set

$$d_w = d_w(G) = \log 5 / \log 2; \quad (2.5)$$

this number governs the space-time scaling of the process  $X$ . Set

$$\Phi(t, x) = t^{-d_f/d_w} \exp(-(x^{d_w} t^{-1})^{1/(d_w - 1)}). \quad (2.6)$$

**Theorem 5** [BP, Theorem 1.5]. *There exists a continuous symmetric function  $p_t(x, y)$ ,  $(t, x, y) \in (0, \infty) \times G \times G$  such that*

(a)  *$p_t$  is the transition density of  $X$  with respect to  $\mu$ :*

$$P_t f(x) = \int f(y) p_t(x, y) \mu(dy), \quad x \in G, \quad f \in bG,$$

(b) *there exist constants  $c_1, \dots, c_4$  such that*

$$c_1 \Phi(t, c_2 |x - y|) \leq p_t(x, y) \leq c_3 \Phi(t, c_4 |x - y|) \quad \text{for } x, y \in G, 0 < t \leq 1, \quad (2.7)$$

(c)  *$u(t, x) = p_t(x_0, x)$  is the fundamental solution of the heat equation on  $G$  with pole at  $x_0$ :*

$$\frac{\partial u}{\partial t} = \mathcal{L}u; \quad u(0, \cdot) = \delta_{x_0}. \quad (2.8)$$

*Remarks.* 1. In [BP] this was actually proved for the process on the unbounded gasket  $\tilde{G} = \bigcup_{n=0}^{\infty} 2^n G$ , and in this case the estimate (2.7) holds for  $t \in (0, \infty)$ .

2. Setting  $x = y$  in (2.7) one obtains

$$p_t(x, x) \approx t^{-d_f/d_w} \quad 0 < t \leq 1. \quad (2.9)$$

(Here  $\approx$  means ‘bounded above and below by constants’.)

Let  $-\lambda_1 \geq -\lambda_2 \geq \dots$  be the eigenvalues of  $\mathcal{L}$ , let  $\varphi_i$  be the corresponding eigenfunctions, and set  $N(\lambda) = \#\{\lambda_i : \lambda_i \leq \lambda\}$ . Then the transition density  $p_t(x, y)$  has an expansion

$$p_t(x, y) = \sum e^{-\lambda_i t} \varphi_i(x) \varphi_i(y);$$

setting  $x = y$  and integrating over  $G$  one obtains for  $t < 1$

$$t^{-d_f/d_w} \approx \int p_t(x, x) \mu(dx) = \sum e^{-\lambda_i t} = \int_0^\infty e^{-\lambda t} N(d\lambda),$$

and inverting the Laplace transform it follows that

$$N(\lambda) \approx \lambda^{d_f/d_w} \quad \text{as } \lambda \rightarrow \infty. \quad (2.10)$$

The quantity  $d_s = 2d_f/d_w$  is referred to as the ‘spectral dimension’ of  $G$ ; from (2.9) and (2.10) one sees that it governs the short-time asymptotics of the transition function, and the frequency of small eigenvalues of  $\mathcal{L}$ .

There is now a well developed ‘machine’ for obtaining estimates on the transition functions of symmetric Markov processes from geometric properties of the underlying space – see [CKS, V] and the papers cited therein. However, this method does not appear to work well in the context of fractal spaces:

- (a) It does not seem to be easy to obtain the correct  $L^1$ -Sobolev inequality as a starting point. In particular, using an isoperimetric inequality does not give the correct value of  $d_s$  – see [O, BBS and BB4].
- (b) Even when one has on-diagonal estimates for  $p_t$ , there are problems in applying E.B. Davies’ method (see [CKS]) for obtaining off-diagonal estimates. One major difficulty is that the measures  $v_{[f,g]}$  associated with the Dirichlet form  $\mathcal{E}(f, g)$  are singular with respect to  $\mu_G$  – see [Ku2].

The proof of Theorem 6(b) proceeds as follows, and it is probable that a similar approach will work for other fractals with  $d_s < 2$ .

(1) Obtain estimates on the potential kernel density  $u_A(x, y)$  for the process  $X$  killed on leaving the region  $A \subset G$ . In particular, one finds that if  $B_m(x) = \{y \in G : |x - y| < 2^{-m}\}$  then

$$u_{B_m(x)}(x, x) \approx (\tfrac{1}{2})^m = (2^{-m})^{d_w - d_f}. \quad (2.11)$$

(2) Obtain estimates on

- (a) the  $P^x$ -law of the times  $T_m$  to leave the region  $B_m(x)$ ,
- (b) the density  $f_{xy}(s)ds = P^x(T_y \in ds)$ .

(3) Combine (1) and (2)(a) to obtain bounds on  $u_\lambda(x, y)$ , the  $\lambda$ -potential density. If  $\lambda = 2^{md_w}$  then the stopping times  $T_m$  and  $R_\lambda$  (an independent negative exponential r.v. with mean  $\lambda^{-1}$ ) are of the same order of magnitude, and therefore so are  $u_\lambda(x, x)$  and  $u_{B_m(x)}(x, x)$ . One has

$$u_\lambda(x, x) \approx \lambda^{\frac{1}{2}d_s - 1}, \quad \lambda \geq 1. \quad (2.12)$$

(4) Applying a Tauberian theorem, and using the fact that since  $X$  is symmetric  $t \mapsto p_t(x, x)$  is decreasing, (2.12) implies that

$$p_t(x, x) \approx t^{-d_s/2}, \quad 0 < t \leq 1. \quad (2.13)$$

(5) The off-diagonal bound follows by using the first entrance decomposition

---


$$p_t(x, y) = \int_0^t f_{xy}(s) p_{t-s}(y, y) ds, \quad (2.14)$$

and substituting the bounds given by (2.13) and (2)(b).

Many properties of the process  $X$  follow from Theorem 6, or the estimates used in the proof. For example, integrating (2.7) over  $G$  one obtains

$$E^x|X_t - x|^2 \approx t^{d_w}, \quad 0 < t < 1, \quad (2.15)$$

while the estimates on the tail of the distribution imply the Lévy Hölder law

$$c \leq \limsup_{\delta \downarrow 0} \sup_{\substack{0 \leq s \leq t \leq 1 \\ |t-s| < \delta}} |X_t - X_s| / h(t-s) \leq c', \quad (2.16)$$

where  $h(u) = u^{1/d_w} (\log(1/u))^{(d_w-1)/d_w}$ . The estimates on  $u_1$  imply that  $X$  has a jointly continuous local time ( $L_t^x, x \in G, t \geq 0$ ) which is Hölder continuous of order  $\frac{1}{2}(d_w - d_f) - \varepsilon$  in the space variable. From this it follows that  $X$  is space-filling:

$$\{X_t, 0 \leq t \leq T\} = G \quad \text{for all sufficiently large } T, \text{ a.s.}$$

I conclude this section with a brief account of further work on the Sierpinski gasket.

A purely analytic approach to the potential theory on  $G$  is given by Metz [M] and Kigami [Ki1]. Kusuoka [Ku2] gives a description of the Dirichlet form  $\mathcal{E}(\cdot, \cdot)$  associated with  $X$  as an infinite matrix product. Two noteworthy consequences of this are, firstly, the result (mentioned above) that the measures  $\nu_{[f,g]}$  are singular with respect to  $\mu$ , and secondly, that the natural filtration of  $X$  is one-dimensional, in the sense of [DV].

Fukushima and Shima [FS], [S] give an explicit description of the eigenvalues  $\lambda_i$  of the operator  $\mathcal{L}$ : a particular consequence is that

$$\lambda^{-d_s/2} N(\lambda) \quad \text{oscillates boundedly as } \lambda \rightarrow \infty. \quad (2.17)$$

Most attention has been given to the Brownian motion on  $G$ . However, a more complete analysis of the analytic behaviour of  $G$  would involve a study of diffusion processes with drift. Kumagai [Kum] has recently studied an interesting class of non-symmetric diffusions on  $G$ , the ‘ $p$ -stream diffusions’.

### 3. Other Finitely Ramified Fractals

It is natural to ask whether the results of the previous section extend to other regular fractals. [G] discusses one other fractal, the Vicsek set, and Krebs [Kr] gave a fairly complete analysis. Lindström [L] defined a large class of self-similar finitely ramified fractals, called ‘nested fractals’, and gave a construction of ‘Brownian motion’ on them. For a full description of this class, see [L]; to

simplify the exposition here I will restrict myself to a subfamily which, however, seems large enough to capture the essential features.

Let  $F_0 = \{a_1, \dots, a_k\}$  be the vertices of a regular  $k$ -sided polygon in  $\mathbb{R}^2$ , and let  $A_0$  be the closed convex hull of  $F_0$ . Let  $\lambda > 1$ ,  $M \geq k$ ,  $a_{k+1}, \dots, a_M$  be points in  $A_0$  and let  $\phi_i$ ,  $1 \leq i \leq M$  be linear contractions with scale factor  $\lambda$  such that

- (a)  $\phi_i(x) = a_i + \lambda^{-1}(x - a_i)$  for  $1 \leq i \leq M$  (3.1)
- (b)  $\phi_i(A_0) \subset A_0$  for  $1 \leq i \leq k$ ,
- (c) the sets  $\phi_i(A_0)$  are disjoint apart from their vertices,
- (d) the set  $A_1 = \bigcup_{i=1}^M \phi_i(A_0)$  is connected,
- (e)  $A_1$  satisfies the same symmetries as  $A_0$ .

The set

$$F = \bigcap_{n=0}^{\infty} \Phi^n(A_0), \quad \text{where } \Phi(\cdot) = \bigcup_{i=1}^M \phi_i(\cdot),$$

is a nested fractal. Note that the condition (c) ensures that  $F$  is finitely ramified. The definition of the  $n$ -complexes and  $n$ -cells is the same as for the Sierpinski gasket; let also  $F_n = \Phi^n(F_0)$ .

An example is the Vicsek set, where  $M = 5$ ,  $\lambda = 3$ ,  $a_1, \dots, a_4$  are the corners of the unit square, and  $a_5 = \frac{1}{3}\sum a_i$  is its centre.

The first major difficulty in constructing a Brownian motion  $X$  on  $F$  is in finding a sequence of decimation invariant random walks. One's first thought might be to take  $Y^m$  to be a nearest neighbour random walk on  $F_m$ , where  $x$  and  $y$  are  $m$ -neighbours if  $x$  and  $y$  are adjacent vertices of a polygon in some  $m$ -cell. However, the example of the Vicsek set shows this will not work: if  $Y^m$  is chosen as above then  $\tilde{Y}^{m-1}$  will make diagonal transitions, while  $Y^{m-1}$  can only make horizontal or vertical ones.

As the random walks  $Y^m$  on  $F_m$  are to be symmetric, it is easiest to define them by choosing a symmetric conductivity matrix  $a_m(x, y)$ ,  $x, y \in F_n$ , and defining the transition matrix  $p_m$  by

$$p_m(x, y) = a_m(x, y)/a_m(x), \quad a_m(x) = \sum_y a_m(x, y). \quad (3.2)$$

Let  $\alpha_1, \dots, \alpha_{k-1}$  satisfy  $\sum \alpha_i = 1$ ,  $\alpha_i \in [0, 1]$ ,  $\alpha_i = \alpha_{k-i}$ , and define

$$a_m(x, y) = \begin{cases} \alpha_i & \text{if } x, y \text{ belong to the same } m\text{-cell and } x \text{ and } y \\ & \text{are } i \text{ steps apart on the circumference;} \\ 0 & \text{otherwise.} \end{cases} \quad (3.3)$$

Let  $P_\alpha^x$  be the law of the random walk  $Y^1$  defined by (3.2) and (3.3) with  $m = 1$ . Let  $T = \min\{r \geq 0 : Y \in F_0 - \{a_1\}\}$ , and set  $f_i(\alpha) = P_\alpha^{a_1}(Y_T = a_i)$ ,  $f(\alpha) = (f_1(\alpha), \dots, f_{k-1}(\alpha))$ . It is not hard to verify that if  $Y^m$  has transition probabilities given by  $\alpha = (\alpha_1, \dots, \alpha_{k-1})$ , then  $\tilde{Y}^{m-1}$  has transition probabilities  $f(\alpha)$ , so that the search for a decimation invariant walk corresponds to the study of the fixed points of the map  $\alpha \rightarrow f(\alpha)$ . Let  $\ell = [k/2]$ .

**Theorem 6** [L]. Set  $K = \{\alpha : \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_\ell\}$ . Then  $f : K \rightarrow K$ , and is continuous. Hence  $f$  has at least one fixed point in  $K$ .

Let  $\alpha \in K$  be a fixed point of  $f$  and  $Y^m$ ,  $m \geq 1$  be the associated random walks. Then there exists a number  $\tau = \tau(\alpha, F)$  such that the processes  $X^m(t) = Y^m(\tau^m t)$  converge weakly to a process  $X$  on  $F$ .

**Theorem 7** [L]. (a)  $X$  is a  $\mu_F$ -symmetric Feller diffusion on  $F$ .

(b)  $X$  is invariant with respect to local isometries of  $F$ .

(c) Let  $\Delta_{\alpha,F}$  be the infinitesimal generator of  $X$ ,  $-\lambda_1 \geq -\lambda_2 \geq \dots$  be the eigenvalues of  $\Delta_{\alpha,F}$  and  $N(\lambda) = \#\{\lambda_i : \lambda_i < \lambda\}$ . Then  $N(\lambda) \approx \lambda^{\frac{1}{2}d_s(\alpha,F)}$ , where

$$d_s(\alpha, F) = 2 \log \tau / \log \lambda. \quad (3.4)$$

*Remarks.* 1. Lindström actually uses nonstandard analysis to construct  $X$ , and the weak convergence is an immediate consequence.

2. If  $\alpha \notin K$  is a fixed point, a limiting process can still be constructed, but it is no longer invariant.

Following Lindström, let us call the numbers  $\lambda, M, \tau$  the *length, mass and time scale factors* for the fractal  $F$ . Then

$$\dim(F) = d_f(F) = \log M / \log \lambda. \quad (3.5)$$

As in the case of the Sierpinski gasket, the number

$$d_w(F) = \log \tau / \log \lambda = 2d_f/d_s \quad (3.6)$$

governs the space/time scaling of  $X$ .

The three ‘dimensions’  $d_f, d_s, d_w$  have been known to mathematical physicists for some time [RT, AO]. Together they appear to summarise fairly completely the behaviour of the process  $X$ . In view of the final relation in (3.6) (which seems to hold in great generality) there are really only two independent quantities,  $d_f(F)$  (called the ‘fractal dimension’) and the spectral dimension  $d_s(F)$ . Of these  $d_f$  is a ‘geometrical’ quantity and for self-similar fractals can be calculated immediately from the length and mass scale factors. On the other hand  $\tau$  and  $d_s$  are ‘analytic’: it appears necessary to solve some equation on  $F$  (or its approximations) to obtain their values.

The relation  $d_s \leq d_f$  (or equivalently  $d_w \geq 2$ ) holds for nested fractals [L], for graphs [Ku1] and for Sierpinski carpets [BB3], and it is probable that it holds in general. No other relation between these dimensions seems likely.

One problem left open by [L] is the uniqueness of the fixed point. That this cannot be altogether straightforward is shown by the case of the Vicsek set, where  $\alpha = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$  and  $\beta = (0, 1, 0)$  are both fixed points. However  $\beta \notin K$  is unstable, and, as remarked above, does not give rise to an invariant process on the fractal. Theorem 7 shows that each fixed point in  $K$  gives rise to a Laplacian  $\Delta_{\alpha,F}$ , which in turn determines an ‘analytic structure’ on  $F$ . The uniqueness of fixed points therefore corresponds to the question as to whether a nested fractal carries only one natural analytic structure.

**Conjecture 8.** *The map  $\alpha \rightarrow f(\alpha)$  has exactly one fixed point in  $K$ , and this fixed point is stable.*

Very recently I have proved this in the case  $k = 4$  for the subfamily of nested fractals described above.

As far as the probabilistic properties of the process  $X$  are concerned, less is known than in the case of the Sierpinski gasket. However, there is little doubt that Theorem 5 holds for many nested fractals, with the single change of  $d_s(F)$  for  $d_s(G)$  in (2.6).

A much wider class of finitely ramified fractals (pcf self-similar sets) has been introduced in [Ki2]. Kigami uses analytic methods, and the question of the existence of a decimation invariant walk is replaced by the existence of a ‘harmonic structure’. Kigami proves that, if a harmonic structure on  $F$  exists, then a Laplacian  $\Delta_F$  and Dirichlet form  $\mathcal{E}_F(\cdot, \cdot)$  can be constructed. However, there are at present neither existence nor uniqueness theorems for harmonic structures on general pcf self-similar sets.

#### 4. Sierpinski Carpets

These are the simplest examples of infinitely ramified fractals, and have been studied in [BB1, BB2, BBS, BB3, and BB4]. I will restrict attention here to the basic Sierpinski carpet (SC) – the methods and results are the same for more general SCs. Let  $H_0$  be the closed convex unit square, and  $H_1 = H_0 \setminus (\frac{1}{3}, \frac{2}{3})^2$ . Repeating this operation (of removing the central square) one obtains a sequence of closed sets  $H_n$  ( $H_n$  consists of  $8^n$  squares each of side  $3^{-n}$ ); the Sierpinski carpet  $H$  is defined by

$$H = \bigcap_{n=0}^{\infty} H_n.$$

It is fairly clear that the methods of Sections 2 and 3 will not work here. Instead, as each  $H_n$  has a connected interior, it is natural to consider continuous approximating processes. Let  $W^n$  be Brownian motion on  $H_n$ , with orthogonal reflection on  $\partial H_n$ , and let  $\alpha_n$  be the maximum mean crossing time of  $H_n$  by  $W^n$ , defined by

$$\alpha_n = \sup_{x \in H_n} E^x(\inf\{t \geq 0 : W^n(t) \in J \cap H_n\}), \quad (4.1)$$

where  $J = \{z = (x, y) \in \mathbb{R}^2 : x = 1 \text{ or } y = 1\}$ . Let  $X_t^n = W^n(\alpha_n t)$ ; the processes  $X^n$  are tight, and so if  $(n_k)$  is a subsequence such that  $X^{n_k}$  converges weakly, the limiting process  $X$  is a continuous  $H$ -valued process. By using tightness of resolvents rather than processes (and, possibly, using a different subsequence), one can ensure that the limiting process is in fact a diffusion.

The length and mass scales of  $H$  are given by  $\lambda_H = 3$ ,  $M_H = 8$ . There exists a number  $\tau$  ( $\tau \approx 10.0118$ ) such that  $\alpha_n \approx (\frac{1}{3}\tau)^n$ :  $\tau$  is the time scale for  $H$ . Let

$$d_f(H) = \frac{\log 8}{\log 3}, \quad d_s(H) = \frac{2 \log 8}{\log \tau};$$

$d_f(H)$  is the Hausdorff dimension of  $H$ , while the following result ensures that  $d_s(H)$  is the spectral dimension of any Brownian motion constructed in this fashion.

**Theorem 9** [BB4]. *Let  $(n_k)$  be any subsequence such that the processes  $X^{n_k}$  converge weakly to a diffusion process  $(X_t, P^x)$ . Then  $X$  is Feller, and  $\mu_H$ -symmetric, and has*

a transition density with respect to  $\mu_H$  which satisfies the conditions (a)–(c) of Theorem 5.

*Remark.* The main gap in the theory for the SC is the lack of any kind of uniqueness result: it is possible (though rather unlikely) that the laws of the  $X^n$  have more than one limit point.

The main steps in the construction of  $X$ , and the proof of Theorem 9, are as follows. Set  $K_n = H_n \cap [0, \frac{1}{2}]^2$ , and let  $L_i$ ,  $1 \leq i \leq 4$  denote the 4 edges of  $H_0$ .

(1) Using a reflection principle, and a path-crossing argument, one can prove the following Harnack inequality [BB1, Theorem 3.1]:

There exists a constant  $c > 0$  such that if  $u : H_n \rightarrow \mathbb{R}_+$  is harmonic in  $H_n - J$  then

$$\sup_{u \in K_n} u(x) \leq c \inf_{u \in K_n} u(x). \quad (4.2)$$

(2) Let

$$R_n^{-1} = \inf \left\{ \int_{H_n} |\nabla u|^2 dx : u = 0 \text{ on } L_0, u = 1 \text{ on } L_2 \right\}; \quad (4.3)$$

thus  $R_n$  is the resistance of  $H_n$  when two opposite edges are short-circuited, and a potential difference is applied between them. Using the Harnack inequality to construct a feasible function for the variational problem one obtains [BB2, Sect. 2]

$$\alpha_n \approx \left(\frac{8}{9}\right)^n R_n. \quad (4.4)$$

(3) Further comparison arguments, of  $H_n$  with associated wire networks, give [BBS, BB3]

$$\frac{1}{4} R_n R_m \leq R_n R_m \leq 4 R_n R_m. \quad (4.5)$$

Standard subadditivity results now imply that there exists a number  $\varrho$  such that  $R_n \approx \varrho^n$ . The time scale  $\tau_H$  is defined by

$$\tau_H = 8\varrho = M_H \varrho. \quad (4.6)$$

(4) The proof of the analogue of Theorem 5 for  $H$  now proceeds in a fashion very similar to the proof in the case of the Sierpinski gasket.

The Equation (4.6), relating the resistance growth of the  $H_n$  with the time and mass scales, is called in the physics literature an Einstein relation, and has been known to physicists for some time – see [RT]. The use of ideas from electrical theory (or, equivalently, from a mathematical viewpoint, the theory of Dirichlet forms) seems to be a powerful tool in this area.

## References

- [AO] Alexander, S., Orbach, R.: Density of states on fractals: “fractons”. *J. Phys. (Paris) Lett.* **43** (1982) L625–L631
- [BB1] Barlow, M.T., Bass, R.F.: Construction of Brownian motion on the Sierpinski carpet. *Ann. Inst. H. Poincaré* **25** (1989) 225–257
- [BB2] Barlow, M.T., Bass, R.F.: Local times for Brownian motion on the Sierpinski carpet. *Prob. Theor. Rel. Fields* **85** (1990) 91–104

- [BB3] Barlow, M.T., Bass, R.F.: On the resistance of the Sierpinski carpet. *Proc. R. Soc. Lond. A.* **431** (1990) 345–360
- [BB4] Barlow, M.T., Bass, R.F.: Transition densities for Brownian motion on the Sierpinski carpet. Preprint 1991
- [BBS] Barlow, M.T., Bass, R.F., Sherwood, J.D.: Resistance and spectral dimension of Sierpinski carpets. *J. Phys. A* **23** (1990) L253–L258
- [BP] Barlow, M.T., Perkins, E.A.: Brownian motion on the Sierpinski gasket. *Prob. Theor. Rel. Fields* **79** (1988) 543–623
- [CKS] Carlen, E.A., Kusuoka, S., Stroock, D.W.: Upper bounds for symmetric Markov transition functions. *Ann. Inst. H. Poincaré* **23** (1987) 245–287
- [DV] Davis, M.H.A., Varaiya, P.: The multiplicity of an increasing family of  $\sigma$ -fields. *Ann. Prob.* **2** (1987) 958–963
- [F] Falconer, K.J.: *Fractal geometry*. Wiley, 1990
- [Fu] Fukushima, M.: *Dirichlet forms and Markov processes*. North-Holland/Kodansha, Tokyo 1980
- [FS] Fukushima, M., Shima, T.: On a spectral analysis for the Sierpinski gasket. Preprint, 1989
- [G] Goldstein, S.: Random walks and diffusion on fractals. In: Kesten, H. (ed.) *Percolation theory and ergodic theory of infinite particle systems* (IMA Math. Appl., vol. 8.) Springer, Berlin Heidelberg New York 1987, pp. 121–129
- [H] Hutchinson, J.E.: Fractals and self-similarity. *Indiana Univ. Math. J.* **30** (1981) 713–747
- [HBA] Havlin, S., Ben-Avraham, D.: Diffusion in disordered media. *Adv. Phys.* **36** (1987) 675–798
- [Ki1] Kigami, J.: A harmonic calculus on the Sierpinski space. *Japan J. Appl. Math.* **6** (1989) 259–290
- [Ki2] Kigami, J.: Harmonic calculus on P.C.F. self-similar sets. Preprint, 1989
- [Kr] Krebs, W.: A diffusion defined on a fractal. PhD thesis, University of California at Berkeley, 1988
- [Ku1] Kusuoka, S.: A diffusion process on a fractal. In: Ito, K., N. Ikeda, N. (eds.) *Symposium on Probabilistic Methods in Mathematical Physics*, Taniguchi, Katata. Academic Press, Amsterdam 1987, pp. 251–274
- [Ku2] Kusuoka, S.: Dirichlet forms on fractals and products of random matrices. *Publ. RIMS Kyoto Univ.* **25** (1989) 659–680
- [Kum] Kumagai, T.: Construction and some properties of a class of non-symmetric diffusion processes on the Sierpinski gasket. Preprint, 1990
- [L] Lindstrøm, T.: Brownian motion on nested fractals. *Mem. Amer. Math. Soc.* **420** (1990)
- [M] Metz, V.: Brownsche Bewegung auf dem Sierpinski gasket aus potentialtheoretischer Sicht. Diplomarbeit, Bielefeld, 1988
- [O] Osada, H.: Isoperimetric dimension and estimates on heat kernels of pre-Sierpinski carpets. *Prob. Theor. Rel. Fields* **86** (1990) 469–490
- [RT] Rammal, R., Toulouse, G.: Random walks on fractal structures and percolation clusters. *J. Phys. Lett.* **44** (1983) L13–L22
- [S] Shima, T.: On eigenvalue problems for the Sierpinski pre-gaskets. To appear *Japan J. Appl. Math.*
- [V] Varopoulos, N.Th.: Isoperimetric inequalities and Markov chains. *J. Funct. Anal.* **63** (1985) 215–239



# Applications of Group Representations to Statistical Problems

*Persi Diaconis*

Department of Mathematics, Harvard University, Science Center, 1 Oxford Street  
Cambridge, MA 02138 USA

## Abstract

Many problems in routine statistical analysis can be interpreted as the decomposition of a representation into irreducible components and the computation and interpretation of the projection of a given vector into these components. Examples include the usual spectral analysis of time series and the statistical analysis of variance. Recently, non-commutative representations have emerged as a practical tool. A variety of approaches have come together to give a unified theory.

## 1. Introduction

The study of a function through the size of its coefficients in an orthogonal expansion is a standard tool. This paper shows that expansions arising from the action of a group on a set occur naturally in a variety of statistical problems.

**Example** (Time Series Analysis). Let  $f(0), f(1), \dots, f(N - 1)$ , be the observed value of a series of events. For example, the  $f(k)$  might be the number of children born in New York City on successive days. Data collected in time often exhibit periodic behavior; New York City birth data looks like this:

$$411, 430, 418, 396, 401, 320, 322, \dots,$$

the pattern of 5 high values followed by two low values persists. This seems surprising until one realizes that about 20% of all births are induced and physicians don't like to work on weekends. Izenman and Zabell (1978) discuss these data.

To find and interpret such periodicities, data  $f(k)$  is often transformed as

$$\hat{f}(j) = \sum_{k=0}^{N-1} e^{2\pi i j k / N} f(k).$$

The data can be recovered by the inversion theorem

$$f(k) = \frac{1}{N} \sum_{j=0}^{N-1} \hat{f}(j) e^{-2\pi i j k / N}.$$

If the transform  $\hat{f}(j)$  is relatively large for only a few values of  $j$ , the inversion theorem shows  $f$  is well approximated by these few periodic components. This gives a simple description of  $f$  and one can try to go back and understand why the few components are large.

This “bump hunting” part of spectral analysis is fully explained in books by Brillinger (1975) and Bloomfield (1976). It is only one part of the story – continuous spectra is the other part (see Tukey (1961) – but it, and its generalizations will dominate the present treatment.

The generalizations presented here involve a finite group  $G$  acting on a finite set  $X$ . Let  $f : X \rightarrow \mathbb{R}$  be a given function. In the example above,  $G$  is a cyclic group of order  $N$  acting on itself by translation. The function  $f(k)$  is the number of children born on day  $k$ . In the example of the next section,  $G$  is the symmetric group acting on itself and  $f(\pi)$  is the number of people in an election who ranked candidates in the permutation  $\pi$ .

Let  $L(X) = \{f : X \rightarrow \mathbb{C}\}$ . This is a vector space on which  $G$  acts by  $sf(x) = f(s^{-1}x)$ . Mashke’s theorem implies that  $L(X)$  splits into a direct sum

$$L(X) = V_0 \oplus V_1 \oplus \cdots \oplus V_J$$

where each subspace is invariant under the group (so  $g \in V_i$  implies  $sg \in V_i$ ) and the pieces are irreducible, so no further splitting is possible. Clearly  $f \in L(X)$  can be written as  $f = \sum_{i=0}^J \hat{f}_i$  with  $\hat{f}_i$  the projection into  $V_i$ .

The empirical finding, to be explored further, is that the subspaces often have simple interpretations and the decomposition of  $f$  into its projection into the  $V_i$  “makes sense”.

**Definition.** *Spectral analysis consists of the computation and interpretation of  $\hat{f}_i$  and the approximation of  $f$  by as few pieces as do a reasonable job.*

This necessarily vague definition encompasses a number of areas of classical statistics. In the next section an example is presented in some detail. Section 3 gives a group-theoretic version of the classical analysis of variance as spectral analysis. Section 4 describes modern work on ANOVA of orthogonal designs as developed by Bailey, Nelder, Speed, Tjur, and their co-workers. That these two approaches lead to the same analysis in nice cases is an important recent result of Rosemary Bailey, Chris Rowley, and their co-workers. This is developed in Section 5. The final section gives pointers to the many topics which couldn’t be covered in this brief review.

Spectral analysis as outlined here is a data analytic variation of ideas suggested earlier by Alan James and Ted Hennen. Hennen’s (1965) monograph is filled with innovative ideas and treats continuous problems as well. Peter Fortini’s (1977) thesis is also an important source of inspiration for the treatment presented here.

Only the rudiments of group representations are needed. The beginning of the books of Ledermann (1987) or Serre (1977) are ample background. I have tried to lay out the background in Diaconis (1988).

## 2. An Example

This section presents data on  $S_5$  the symmetric group on 5 letters. The data arise from an election of the American Psychological Association. This organization asks its membership to rank order 5 candidates for president. Here  $G = S_5$  and  $f(\pi)$  is the number of voters choosing rank order  $\pi$ . For example,  $f\begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 5 & 4 & 3 & 2 & 1 \end{pmatrix} = 29$  so 29 voters ranked candidate one 5th, candidate two 4th and so on. The data is shown in Table 1.

Let  $\varrho : S_5 \rightarrow GL_5(\mathbb{R})$  be the usual 5-dimensional permutation representation. Thus  $\varrho(\pi)$  is a  $5 \times 5$  matrix with  $(i, j)$  entry 1 if  $\pi(i) = j$  and zero otherwise. The Fourier transform of  $f$  at  $\varrho$  is the matrix

**Table 1.** American Psychological Association election data

Ranking	No. of votes cast of this type	Ranking	No. of votes cast of this type	Ranking	No. of votes cast of this type	Ranking	No. of votes cast of this type
54321	29	43521	91	32541	41	21543	36
54312	67	43512	84	32514	64	21534	42
54231	37	43251	30	32451	34	21453	24
54213	24	43215	35	32415	75	21435	26
54132	43	43152	38	32154	82	21354	30
54123	28	43125	35	32145	74	21345	40
53421	57	42531	58	31542	30	15432	40
53412	49	42513	66	31524	34	15423	35
53241	22	42351	24	31452	40	15342	36
53214	22	42315	51	31425	42	15324	17
53142	34	42153	52	31254	30	15243	70
53124	26	42135	40	31245	34	15234	50
52431	54	41532	50	25431	35	14532	52
52413	44	41523	45	25413	34	14523	48
52341	26	41352	31	25341	40	14352	51
52314	24	41325	23	25314	21	14325	24
52143	35	41253	22	25143	106	14253	70
52134	50	41235	16	25134	79	14235	45
51432	50	35421	71	24531	63	13542	35
51423	46	35412	61	24513	53	13524	28
51342	25	35241	41	24351	44	13452	37
51324	19	35214	27	24315	28	13425	35
51243	11	35142	45	24153	162	13254	95
51234	29	35124	36	24135	96	13245	102
45321	31	34521	107	23541	45	12543	34
45312	54	34512	133	23514	52	12534	35
45231	34	34251	62	23451	53	12453	29
45213	24	34215	28	23415	52	12435	27
45132	38	34152	87	23154	186	12354	28
45123	30	34125	35	23145	172	12345	30

$$\hat{f}(\varrho) = \sum_{\pi} \varrho(\pi) f(\pi).$$

This has  $(i, j)$  entry the number of people ranking candidate  $i$  in position  $j$ . This natural summary is shown in Table 2 where entries are divided by the total number of voters to give proportions.

**Table 2.** Percentage of voters ranking candidate  $i$  in position  $j$

Candidate	Rank				
	1	2	3	4	5
1	18	26	23	17	15
2	14	19	25	24	18
3	28	17	14	18	23
4	20	17	19	20	23
5	20	21	20	19	20

The largest number 28 in the  $(3, 1)$  position shows 28 percent of the voters ranked candidate 3 first. Candidate 3 also had some “hate vote”: 23 percent ranked 3 last. This first order summary is the first thing anyone analyzing such data would try. It is natural to ask if it captures the essence of the voting pattern or if there is more to be learned.

The data is summarized by  $f \in L(S_5)$ . This last vector space splits into 7 invariant subspaces in its isotypic decomposition shown in Table 3.

**Table 3.** Decomposition of the regular representation

$M =$	$V_1$	$\oplus$	$V_2$	$\oplus$	$V_3$	$\oplus$	$V_4$	$\oplus$	$V_5$	$\oplus$	$V_6$	$\oplus$	$V_7$
Dim 120	1		16		25		36		25		16		1
SS/120	2286		298		459		78		27		7		0

Table 2 amounts to looking at the projection of  $f$  into  $V_1 \oplus V_2$ . If  $L(S)$  is treated as an inner product space with  $\langle f | g \rangle = \sum f(\pi)g(\pi)$  the function  $f$  decomposes into the pieces of its orthogonal projection. The norm square of  $f$  decomposes into the norm squared of its projections by Pythagoras’s theorem. These squared lengths are shown in the last line of Table 3. As usual the largest contribution comes from the projection onto the constant functions. There is also a large projection onto the space  $V_3$ . This projection is not captured in the summary of Table 2.

The space  $V_3$  is made up of “2nd order functions”, a typical such being  $\pi \mapsto \delta_{\{j, j'\}} \{\pi(i), \pi(i')\}$  which is 1 if the unordered pair  $\{\pi(i), \pi(i')\} = \{j, j'\}$ . The span of all 2nd order functions, orthogonal to  $V_1 \oplus V_2$  make up  $V_3$ . In group representations language,  $V_3$  is the isotypic subspace corresponding to the partition 3, 2. This  $V_3$  is 25-dimensional. To understand the projection of  $f$  into  $V_3$  a device of Colin Mallows was used. The function  $f$  corresponding to the data is projected onto  $V_3$ .

The inner product of this projection with functions  $\delta_{\{j,j'\}}\{\pi(i), \pi(i')\}$  is then reported. The pairs  $\{i, i'\}, \{j, j'\}$  can be chosen in 10 ways each. The 100 inner products are shown in Table 4.

**Table 4.** Second order, unordered effects

Candidate	Rank									
	1, 2	1, 3	1, 4	1, 5	2, 3	2, 4	2, 5	3, 4	3, 5	4, 5
1, 2	-137	-20	18	140	111	22	4	6	-97	-46
1, 3	476	-88	-179	-209	-147	-169	-160	107	128	241
1, 4	-189	51	113	24	-9	98	99	-65	23	-146
1, 5	-150	57	47	45	43	49	56	-48	-53	-48
2, 3	-42	84	19	-61	30	-16	82	-76	-39	72
2, 4	157	-20	-43	-25	-93	-76	-56	8	38	112
2, 5	22	-44	7	15	-117	69	25	62	99	-138
3, 4	-265	-7	72	199	39	140	85	19	-52	-233
3, 5	-169	10	88	70	78	44	47	-51	-36	-80
4, 5	296	-24	-142	-130	-5	-163	-128	38	-9	267

To explain, consider the  $\{1, 3\}\{1, 2\}$  entry 476. This is the largest number in the table. It means that there is a large positive effect for ranking candidates one and three in the first two positions of the ballot. The last entry in row  $\{1, 3\}$  shows that these candidates also had a lot of hate vote.

A very similar picture occurs for the last row of the table. With these observations, the tables pattern becomes apparent. The American Psychological Association consists of two groups, academicians and clinicians who are on uneasy terms. Candidates  $\{1, 3\}$  are from one group,  $\{4, 5\}$  from the other. Very few voters cross ranks so there is a large negative effect for ranking, say  $\{1, 4\}$ , first and second (or fourth and fifth). These observations are the main structure not revealed by the first order analysis.

In studying data as we have above it is natural to ask if the data were collected again, would the same patterns arise. I will not go into the details here, but a variety of stochastic analyses suggest that the natural scale of variability in Table 4 is  $\pm 50$ , so the patterns observed are believable.

Further details of this analysis are given in Diaconis (1989). Diaconis and Smith (1989) give a different set of applications for these group-theoretic decompositions.

### 3. Analysis of Variance (ANOVA)

Consider data cross classified by  $I$  levels of one variable and  $J$  levels of a second variable. The observed data is then a function  $f(i, j)$  from  $X = \{(i, j) : 1 \leq i \leq I, 1 \leq j \leq J\}$  into  $\mathbb{R}$ . The product  $S_I \times S_J$  acts on  $X$  and  $L(X)$  splits into

$$L(X) = V_0 \oplus V_1 \oplus V_2 \oplus V_3$$

$$\dim I \times J - 1 \quad I - 1 \quad J - 1 \quad (I - 1)(J - 1)$$

where  $V_0$  is the space of constant functions,  $V_1$  is the space of row functions  $f(i, j) = f(i, j')$ ,  $V_2$  is the space of column functions, and  $V_3$  is the space orthogonal to  $V_0 \oplus V_1 \oplus V_2$ . The projection of  $f$  onto  $V_0 \oplus V_1 \oplus V_2$  can be interpreted as the least squares approximation to  $f$  of form  $f(i, j) = \alpha + \beta_i + \gamma_j$  for constants  $\alpha, \beta_i, \gamma_j$ .

This, and many more complex variants are known collectively in statistical literature as the Analysis of Variance. The classical book by Sheffe (1959) is still the best treatment of this widely used subject.

The group-theoretic treatment of ANOVA was pioneered by Alan James and Ted Hennen with important later work by Peter Fortini. Group theory is useful in analyzing more complex designs where the appropriate decomposition is not so easy to guess at. The dimensions and projections of various subspaces can be computed by character theory. Diaconis (1988) reviews these topics.

Even in the simple example given above, thinking group-theoretically has something to offer: Instead of  $S_I \times S_J$  one can consider  $S_I \times C_J$  or  $C_I \times C_J$ , with  $C_I$  a cyclic group of order  $I$ . These groups act transitively and their use would be appropriate if the order of the corresponding rows or columns matters. For example, if the rows of the table were birds, and the columns months of the year, with  $(i, j)$  entry the number of birds of type  $i$  cited in month  $j$  (all in a given location) the decomposition by  $S_I \times C_{12}$  would be appropriate.

Carrying out the projections involves calculating the Fourier transform at many different irreducible representations. In Diaconis and Rockmore (1990) a non-commutative analog of the fast Fourier transform has been developed and used to make these computations efficient. Similar work is being developed by Beth (1984) and Clausen (1989a, b). Historically, the FFT on  $C_2^d$  was first developed by Yates (1937) to analyze multi way tables.

## 4. Modern ANOVA

Analysis of variance has developed along non group-theoretic lines. In this section a survey of works by Rosemary Bailey, Chris Rowley, John Nelder, Tue Tjur, Terry Speed, and their co-workers is given. The next section shows how the present treatment and group-theoretic treatment interact.

Begin with a finite set  $X$ . Let  $L(X)$  be the real-valued functions on  $X$ . A design  $\mathcal{D}$  is a set of partitions  $F$  of  $X$ . For example, in ANOVA with repeated observations in a cell  $X = \{(i, j, k) : 1 \leq i \leq I, 1 \leq j \leq J, 1 \leq k \leq n_{ij}\}$  the design might be taken as  $\mathcal{D} = \{U, R, C, R \wedge C, E\}$  where  $U$  is the universal partition with one block.  $R$  is the row partition  $(i, j, k) \xrightarrow{R} (i', j', k')$  iff  $i = i'$ ,  $C$  is the column partition,  $R \wedge C$ , the minimum of  $R$  and  $C$ , has indices equivalent if they are in the same row and column, and  $E$  is the partition into singletons.

For each partition  $F \in \mathcal{D}$ , let  $L_F = \{f \in L(X) \text{ which are } F \text{ measurable}\}$ . The projection  $P_F : L(X) \rightarrow L_F$  is defined by the averaging matrix

$$(P_F)_{xy} = \begin{cases} 1/|f| & x, y \in f \text{ with } f \text{ a block in } F \\ 0 & \text{otherwise.} \end{cases}$$

Two partitions,  $F, G \in \mathcal{D}$  are called orthogonal if their subspaces are geometrically orthogonal, or if  $P_F P_G = P_G P_F$ . For the ANOVA example,  $R$  and  $C$  are orthogonal provided  $n_{ij} = |X|n_{i+}n_{+j}$  with  $n_{i+} = \sum_j n_{ij}$ . An orthogonal design has all factors orthogonal.

In recent years, orthogonal designs with the set of factors closed under maxima have come to be seen as a useful class with a unified theory. Adding maxima of orthogonal factors preserves orthogonality, so a design can always be completed in this way. Such designs are called Tjur designs because of the following basic result.

**Theorem** (Tjur (1984)). *A Tjur design admits a unique decomposition*

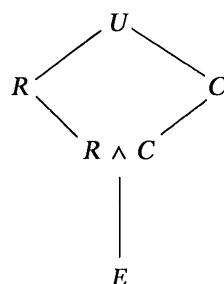
$$L(X) = \bigoplus_{G \in \mathcal{D}} V_G$$

with the property that

$$L_F = \bigoplus_{\substack{G \in \mathcal{D} \\ G \geq F}} V_G.$$

The projection of a given  $f \in L(X)$  onto the various  $V_G$  constitute the analysis of variance for a Tjur design. The point of the decomposition is this: it is easy to complete the projection  $P_F$  of  $f$  onto  $L_F$ . The partially ordered set of factors then allows the computation onto the  $V_G$  by subtraction which amounts to Möbius inversion in the poset.

In the basic ANOVA example, the factors can be diagrammed as:



Given  $f \in L(X)$  one computes the projection onto the constant functions by  $P_U$ . The projection onto the row effects space  $V_R$  is given by  $P_R - P_U$ . The projection onto the column effects space  $V_C$  is given by  $P_C - P_U$ . The projection onto the residual or interaction space  $V_{R \wedge C}$  is given by  $P_{R \wedge C} - P_R - P_C + P_U$ .

More generally, the projection onto  $V_G$  is given by  $\sum_{F \geq G} \mu(F, G)P_F$  with  $\mu$  the Möbius function of the partially ordered set of factors.

This is an easy algorithm which is used by many large computer programs (e.g., Genstat) to analyze designed experiments. A splendid treatment of this point of view appears in Tjur (1984).

## 5. Groups and Modern ANOVA

Sections 3 and 4 above present two different approaches to the analysis of designed experiments. In both, data are represented as  $f \in L(X)$ . In the group case, a group  $G$  is found acting on  $X$  and the analysis consists of decomposing  $L(X)$  and computing projections. In the second case, a collection of partitions is produced and one uses the splitting of  $L(X)$  into parts indexed by these partitions.

It is natural to ask about the relation between these approaches. This problem has been solved for a large class of examples in recent work of Bailey, Prager, Rowley, and Speed (1983). Their work intertwines the two approaches. It is also important group-theoretically in providing examples where the Fourier transform can be computed using the simple averaging and difference algorithm outlined before.

To describe their result, let  $X$  be a finite set and  $\mathcal{D} = \{F\}$  a design, or set of partitions of  $X$ . Assume that the blocks of each  $F \in \mathcal{D}$  have uniform size, that all partitions  $F, G$  in  $\mathcal{D}$  are orthogonal, that  $\mathcal{D}$  is closed under max and min, and finally that  $\mathcal{D}$  forms a distributive lattice under max and min.

These assumptions include many complex classical cases. However, adding the minimum of two partitions to an orthogonal design can destroy orthogonality.

As automorphism of a design  $\mathcal{D}$  is a 1-1 map  $\pi : X \rightarrow X$  such that for each  $F \in \mathcal{D}$ , if  $x$  and  $y$  are in the same block of  $F$  then  $\pi(x)$  and  $\pi(y)$  are in the same block of  $F$ . The set of all automorphisms of  $\mathcal{D}$  is called the automorphism group of  $\mathcal{D}$ .

Bailey, Prager, Rowley, and Speed (1983) did three main things:

I) They identified the automorphism group of  $\mathcal{D}$  as a generalized wreath product of symmetric groups indexed by a partially ordered set. These generalized wreath products have been extensively developed because of their role in the algebraic theory of semi-groups (Krohn-Rhodes theory). A marvelous introduction to this theory appears in Wells (1976). The result also builds on previous work by Silcock (1977).

II) They identify the characters of the automorphism group that appear in the representation  $L(X)$ .

III) They show that the group-theoretic and partition based analysis agree.

Much further work is not reported here. For example, they determine the commuting algebra of  $L(X)$ , give a natural language for describing the groups and decomposition, and finally they make the link with the large body of statistical work in a useful way.

Each approach has problems where it seems to be the superior mode of analysis. The approach by partitions works for some designs without enough symmetry to permit a useful group-theoretic analysis. For example, consider a 2-way array with  $n_{ij}$  entries per cell where  $n_{ij}$  are as shown:

1	2	3
2	4	6
3	6	9

Here it does not make sense to permute the rows or columns. However, the design is orthogonal and permits a straightforward analysis.

In the other direction, block designs are a widely used class of design which are not orthogonal. As an example, consider an experiment in which  $v$  levels of vanilla are to be compared to help decide how much to put into ice cream. If one asks people to taste many ice cream cones, they all taste like colored sugar water. Thus suppose people are asked to taste  $k < v$  flavors. A complete block design involves  $\binom{v}{k}$  people each of whom tastes  $k$  levels of vanilla. Suppose the response is a rating between 0 and 100. This yields  $k \binom{v}{k}$  responses in total. The underlying set

$X = \{(i, s) : 1 \leq i \leq v, |s| = k, i \in s\}$  the responses give  $f: X \rightarrow \mathbb{R}$ .

The two natural partitions are for treatments and blocks. Thus  $(i, s) \xsim{T} (i', s')$  if  $i = i'$  and  $(i, s) \xsim{B} (i', s')$  if  $s = s'$ . These two partitions do not yield an orthogonal design. Indeed, here  $T \wedge B = U$  and the condition for orthogonality can be stated as  $n_{is}|X| = n_i n_s$  where

- $n_i$  is the number of elements in  $X$  receiving treatment  $i$  (so  $n_i = \binom{v-1}{k-1}$ )
- $n_s$  is the number of elements in  $X$  in block  $s$  (so  $n_s = k$ )
- $n_{is}$  is the number of elements in  $X$  with  $x = (i, s)$  (so  $n_s = 1$  if  $i \in s$ , 0 if  $i \notin s$ ).

It follows that  $n_{is}|X| \neq n_i \cdot n_s$  for  $i \notin s$ .

The group-theoretic analysis of this kind of data is straightforward but picks up aspects not developed in earlier analysis. The automorphism group can be identified with  $S_v$ . The representation  $L(X)$  decomposes into a treatment space and a block space, but it also includes new pieces which may be interpreted as the effect of taster's rating by comparison. Fortini (1977) or Diaconis (1988) give further detail.

## 6. Other Topics

This section gives pointers to closely related research which cannot be adequately covered due to space limitations.

### 6.1 Stochastic Models

The approach to spectral analysis outlined above begins with data and a group. Almost all of the statistical literature begins with a probability model and presents the projections as estimates of parameters in a model. For example, for two way analysis of variance with one observation per cell the model would be written as

$$f(i, j) = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$$

where  $\mu, \alpha_i, \beta_j$  are parameters to be estimated (and  $\sum \alpha_i = \sum \beta_j = 0$  to yield identifiable parameters). The  $\varepsilon_{ij}$  are errors, or disturbance terms, which are usually assumed to be independent random variables with mean 0 and constant variance. The least squares estimates of these parameters are the projections described earlier.

Assuming a model leads to well understood ways to quantify standard errors for the estimates. It also allows analysis of data with no symmetry at all. Further, if more careful specifications are made on the distribution of the error terms, a variety of other estimates of the parameters become available.

One of the nice results of the past 10 years is a complete understanding of all possible covariance structures for the error terms which lead to the original projections being efficient estimators. This, and the closely related subject of general balance are treated by Speed (1987), or Bailey and Rowley (1990).

Rosemary Bailey has developed a more elaborate theory which allows incorporation of the randomization aspect of many designed experiments. Her treatment provides separate provision for treatment and design aspects. The theory makes extensive use of group theory and is well related with statistical practice. A recent survey with extensive pointers to other work is Bailey (1990).

There has also been an extensive development which assumes that the errors are Gaussian. The leading work here comes out of the Danish school. Andersson and Perlman (1989) is an important paper with pointers to other work.

## 6.2 Bayesian Methods and Shrinkage Estimators

Once a model is specified, the Bayesian approach to statistics proceeds by putting a prior distribution on the parameters and then using observations to get a posterior distribution. There has been very little work on analyzing the kind of data discussed above from a Bayesian perspective. Dawid (1988) presents some results as do Box and Tiao (1973) and Diaconis, Eaton, and Lauritzen (1991) but much remains to be done.

One of the exciting findings of recent statistical research has been the understanding that when many parameters must be estimated, the classical projection estimates can be uniformly improved. The improvement depends on the assumption of a model. Our current understanding of a reasonable way to go after the improvement involves a Bayesian (or empirical Bayesian) treatment of the problem. Again there has not been much work on shrinkage estimates for designed experiments but Bock (1975) and George (1986) are good starts.

## 6.3 Messy Data

Real data often contains a few stray or wild values that will foul up the classical linear estimates. There is a growing theory of robust statistics surveyed in Huber (1981) or Hoaglin, Mosteller and Tukey (1983, 1985). Much remains to be done in specializing the results available for robust regression to the demands of a complex designed experiment.

Tukey (1977) has developed robust analyses in much the same data analytic spirit as presented here. He supplements these with an extensive residual analysis for ferreting out non-linearities and wild values. He also gives techniques for fitting non-linear models such as

$$f(i, j) = \alpha + \beta_i + \gamma_j + \delta\beta_i\gamma_j + \varepsilon_{ij}.$$

There is active work on non-linear substitutes for classical least squares estimates. Projection pursuit, as developed by Friedman, Stutzle, and Schroeder (1984) or Huber (1985) is one of many varieties. Again, the adaption to analysis of designed experiments is largely open.

Finally, missing data is an annoying part of real statistical analysis. A neat design can become a nightmare with symmetry destroyed. The E.M. Algorithm is now a standard tool for beginning to deal with this problem. Dempster, Laird, and Rubin (1977) is a good reference and Little and Rubin (1987) is a comprehensive guide to the state of the art.

## 6.4 Final Words

The problem with a model based analysis is that usually the model is simply made up out of whole cloth, from linearity through stochastic assumptions. While in principle assumptions can be checked, in my experience they are wildly misused. See Freedman (1986, 1987) for an extensive discussion.

Complex models have the further disadvantage that their parameters are not simple-to-interpret averages but estimates of rather complex quantities whose interpretation depends crucially on the correctness of the model. The spectral estimate proposed in the first sections of this paper are relatively easy to understand averages.

Finally, as for block designs, group-theoretic considerations can lead to different analyses and new models. There is clearly much to be done in combining the best features of the various theories and then confronting them with reality.

## References

- Andersson, S., Perlman, M. (1988): Lattice models for conditional independence in a multivariate normal distribution. Technical report 155, Department of Statistics, University of Washington
- Bailey, R. (1990): A unified approach to design of experiments. *J. Roy. Statist. Soc. A* **144**, 214–223
- Bailey, R., Praeger, T., Rowley, C., Speed, T. (1983): Generalized wreath products of permutation groups. *Proc. London Math. Soc.* **47**
- Bailey, R., Rowley, C. (1990): General Balance and Treatment Permutations. *Lin. Alg. Appl.* **127**, 183–225
- Beth, T. (1984): *Versfahren der Schnellen Fourier-Transform*. Teubner, Stuttgart
- Bloomfield, P. (1976): Fourier analysis of time series. An introduction. Wiley, New York
- Bock, M.E. (1975): Minimax estimators of the mean of a multivariate normal distribution. *Ann. Statist.* **3**, 209–218
- Box, G., Tiao, G. (1973): Bayesian inference in statistical analysis. Addison-Wesley, Reading, Mass
- Brillinger, D. (1975): Time series, data analysis and theory. Holt, Rinehart, and Winston, New York
- Clausen, M. (1989a): Fast Fourier transforms for meta-Abelian groups. *SIAM J. Comput.* **18**, 584–593
- Clausen, M. (1989b): Fast generalized Fourier transforms. *J. Theoret. Comput. Sci.* **67**, 55–63
- Dawid, P. (1988): Symmetry models and hypotheses for structured data layouts. *J. Roy. Statist. Soc. B* **50**, 1–26

- Dempster, A., Laird, N., Rubin, D. (1977): Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc. B* **39**, 1–38
- Diaconis, P. (1988): Group representations in probability and statistics. Institute of Math. Statistics, Hayward, CA
- Diaconis, P. (1989): A generalization of spectral analysis with application to ranked data. *Ann. Statist.* **17**, 949–979
- Diaconis, P. and Smith, L. (1989): Residual analysis for discrete longitudinal data. Technical report. Department of Statistics, Stanford University
- Diaconis, P., Rockmore, D. (1990): Efficient computation of the Fourier Transform on finite groups. *J. Amer. Math. Soc.* **3**, 297–332
- Diaconis, P., Eaton, M., Lauritzen, S. (1991): Finite de Finetti theorems in linear models and multivariate analysis. Technical report Dept. of Mathematics, Aalborg University. To appear in *Scand. J. Statist.*
- Fortini, P. (1977): Representation of groups and the analysis of variance. Ph.D. Thesis, Department of Statistics, Harvard University
- Friedman, J.H., Schroeder, A., Stutzle, W. (1984): Projection pursuit density estimation. *J. Amer. Statist. Assoc.* **79**, 599–608
- Freedman, D. (1987): As others see us: A case study in path analysis. *J. Educ. Statist.* **12**, 101–206
- Freedman, D., Navidi, W. (1986): Regression models for adjusting the 1980 census. *Statist. Sci.* **1**, 3–39
- George, E. (1986): Minimax multiple shrinkage estimation. *Ann. Statist.* **14**, 188–205
- Hannen, E.J. (1965): Group representations and applied probability. Methuen, New York
- Hoaglin, D., Mosteller, F., Tukey, J. (1983): Understanding robust and exploratory analysis. Wiley, New York
- Hoaglin, D., Mosteller, F., Tukey, J. (1985): Exploring data: Tables, trends, and shapes. Wiley, New York
- Huber, P. (1981): Robust statistics, Wiley, New York
- Huber, P. (1985): Projection pursuit. *Ann. Statist.* **13**, 435–525
- Izenman, A., Zabell, S. (1978): Babies and the blackout: the genesis of a misconception. Technical report 38, Dept. of Statistics, University of Chicago
- James, A. (1957): The relationship algebra of an experimental design. *Ann. Math. Statist.* **27**, 993–1002
- James, A. (1982): Analysis of variance determined by symmetry and combinatorial properties of zonal polynomials. In G. Kallianpur et al. (eds.) *Statistics and probability: essays in honor of C.R. Rao*. North-Holland, New York
- Ledermann, W. (1987): Introduction to group characters (2nd ed.). Cambridge University Press, Cambridge
- Little, R., Rubin, D. (1987): Statistical analysis with missing data. Wiley, New York
- Scheffe, H. (1959): Analysis of variance. Wesely, New York
- Serre, J.P. (1977): Linear representations of finite groups. Springer, Berlin Heidelberg New York
- Silcock, H.L. (1977): Generalized wreath products and the lattice of normal subgroups of a group. *Algebra Universalis* **7**, 361–372
- Speed, T. (1987): What is an analysis of variance? *Ann. Statist.* **15**, 885–941
- Tjur, T. (1984): Analysis of variance models in orthogonal designs. *Int. Statist.* **52**, 33–82
- Tukey, J. (1961): Discussion, emphasizing the connection between analysis of variance and spectrum analysis. *Technometrics* **3**, 191–219
- Tukey, J. (1977): Exploratory data analysis. Addison-Wesley, Reading, Mass
- Wells, C. (1976): Some applications of the wreath product construction. *Amer. Math. Monthly* **83**, 317–338
- Yates, F. (1937): The design and analysis of factorial experiments. Imperial Bureau of Soil Science, Harpenden, England

# Stochastic Models of Growth and Competition

Richard Durrett

Department of Mathematics, White Hall, Cornell University, Ithaca, NY 14853, USA

In this paper we will describe recent results on four interacting particle systems that model the growth and competition of plant species or the spread of an epidemic or forest fire. In each system there is a collection of sites, the  $d$ -dimensional integer lattice, that at each time  $t \in [0, \infty)$  can be in one of a finite number of states, so the state of the process at time  $t$  is a function  $\xi_t : \mathbb{Z}^d \rightarrow \{0, 1, \dots, k\}$ . The time evolution is described by declaring that each site changes its state at a rate that depends upon the states of a finite number of neighboring sites. Here, we say that something happens at rate  $r$  if the probability of an occurrence between times  $t$  and  $t + h$  is  $rh + o(h)$ .

## 1. The Basic Contact Process

In this model  $\xi_t : \mathbb{Z}^d \rightarrow \{0, 1\}$ , we think of 0 as vacant and 1 as occupied by a “particle,” and the system evolves as follows:

- (i) Particles die at rate one, give birth at rate  $\beta$ .
- (ii) A particle born at  $x$  is sent to a  $y$  chosen at random from the  $2d$  nearest neighbors  $\{y : \|x - y\|_1 = 1\}$ .
- (iii) If  $y$  is occupied then the birth is suppressed.

Rule (iii) says that there can be at most one particle per site. This is a reasonable constraint if you are thinking of the spread of a plant species but this realism makes the model very difficult to analyze. Let  $\xi_t^A$  be the state at time  $t$  when initially  $\xi_0^A(x) = 1$  if and only if  $x \in A$ , and let  $\tau^A = \inf\{t : \xi_t^A \equiv 0\}$ . If there are no particles then none can be born, so  $\xi_t^A \equiv 0$  for all  $t \geq \tau^A$ . In words, the “all 0” state is an *absorbing state* and we say the system *dies out* at time  $\tau^A$ .

The first question to be addressed is “When does the system have positive probability of not dying out starting from a single occupied site?” or “When is  $P(\tau^{\{0\}} = \infty) > 0$ ?”. It suffices to use a single occupied site as an initial configuration since  $P(\tau^{\{0\}} = \infty) = 0$  implies  $P(\tau^A = \infty) = 0$  for all finite  $A$ . Now, increasing  $\beta$  improves the chances for survival, so it should be clear that there is a critical value

$$\beta_c = \inf\{\beta : P(\xi_t^0 \neq 0 \text{ for all } t) > 0\}.$$

If we delete rule (iii) from the definition, the resulting system is called a branching random walk and has  $\beta_c = 1$ . That is, in order for a branching random walk to survive it is sufficient to have a birth rate larger than the death rate. Since in the contact process some of the birth rate will be wasted on occupied sites, this proves the easy half of the following result.

**Theorem 1A.**  $1 < \beta_c(\mathbb{Z}^d) \leq 4$ .

The lower bound is due to Harris (1974), the upper bound to Holley and Liggett (1978). Both bounds are reasonably accurate. Numerical results (see Brower, Furman, and Moshe (1978)) suggest that  $\beta_c(\mathbb{Z}) \approx 3.299$  and  $\beta_c(\mathbb{Z}^2) \approx 1.645$ , and it has been shown (see Holley and Liggett (1981) or Griffeath (1983)) that  $\beta_c(\mathbb{Z}^d) \rightarrow 1$  as  $d \rightarrow \infty$ .

Once it was established that  $\beta_c \in (0, \infty)$ , attention turned to “What does the process look like when it does not die out?” To answer this question we begin by introducing a special property of the contact process called *duality*

$$P(\xi_t^A(x) = 0 \text{ for } x \in B) = P(\xi_t^B(x) = 0 \text{ for } x \in A).$$

An immediate consequence of duality is that if we start from  $\xi_0^1(x) \equiv 1$  then  $\xi_t^1 \Rightarrow \xi_\infty^1$ . Here,  $\Rightarrow$  is short for *converges weakly* and means that

$$P(\xi_t^1(x) = 0 \text{ for } x \in B) \rightarrow P(\xi_\infty^1(x) = 0 \text{ for } x \in B)$$

for all finite sets B. To prove the weak convergence we set  $A = \mathbb{Z}^d$  in the duality equation to get

$$P(\xi_t^1(x) = 0 \text{ for } x \in B) = P(\xi_t^B(x) \equiv 0)$$

which increases to a limit as  $t \rightarrow \infty$ , since “all 0” is an absorbing state. It follows from standard results (see Chapter 1 of Liggett (1985)) that  $\xi_\infty^1$  is a stationary distribution for the contact process, i.e., if we start the process with this distribution it has this distribution for all time.

At the other extreme, the point mass on the “all 0” state,  $\delta_0$ , is a trivial stationary distribution. Letting  $B = \{y\}$  and  $t \rightarrow \infty$  in the duality relation gives

$$P(\xi_\infty^1(y) = 0) = P(\tau_t^{(y)} < \infty),$$

so  $\xi_\infty^1 = \delta_0$  if the contact process dies out, but is a nontrivial stationary distribution if the contact process survives. The next result, called the *complete convergence theorem* implies that  $\xi_\infty^1$  is the only nontrivial stationary distribution.

**Theorem 1B.**  $\xi_t^A \Rightarrow P(\tau^A < \infty) \delta_0 + P(\tau^A = \infty) \xi_\infty^1$ .

In words, when the process dies out it looks dead, but when it survives and t is large it looks like the system starting from all sites occupied.

The last result took fifteen years to evolve to its current form. Harris (1974), Griffeath (1978), Durrett (1980), Durrett and Griffeath (1982), and Durrett and

Schonmann (1987) proved increasingly more general results before Bezuidenhout and Grimmett (1990) finished the problem and in addition proved

**Theorem 1C.** *When  $\beta = \beta_c$ ,  $P(\tau^{\{0\}} = \infty) = 0$ .*

In words, the contact process dies out at the critical value. For applications (including some we will make in this paper) it is worthwhile to note that all the results in this section hold if (ii) is replaced by

(ii) A particle born at  $x$  is sent to a  $y$  chosen at random from  $x + \mathcal{N}$ .

if we assume  $\mathcal{N}$  is (a) *symmetric* with respect to reflection in any coordinate plane, and (b) *irreducible*, i.e., the group generated by  $\mathcal{N}$  is  $\mathbb{Z}^d$ .

## 2. Multitype Contact Processes

It is well known, even to mathematicians, that there is more than one type of plant, so it is natural to generalize the contact process to have two (or more) types of particles. In this model, the state at time  $t$   $\xi_t : \mathbb{Z}^d \rightarrow \{0, 1, 2\}$  and we think of 0 as vacant and 1 and 2 as occupied by pine and maple trees respectively. With this in mind we formulate the evolution as follows:

- (i) Particles of type  $i$  die at rate one, give birth at rate  $\beta_i$ .
- (ii) A particle born at  $x$  is sent to a  $y$  chosen at random from  $x + \mathcal{N}$  where  $\mathcal{N}$  is symmetric and irreducible.
- (iii) If  $y$  is occupied then the birth is suppressed.

When only one type of particle is present the system reduces to the basic contact process so if  $\beta_1, \beta_2 > \beta_c(\mathbb{Z}^d)$  then there are three trivial equilibria:  $\delta_0$ ,  $\mu_1$  and  $\mu_2$ , where  $\mu_i$  is the limit starting from  $\xi_t(x) \equiv i$ . The main question to be answered about the new system is: "Is there a nontrivial stationary distribution?", i.e., one that concentrates on configurations that contain both 1's and 2's. The first result is a negative one.

**Theorem 2A.** *If  $\beta_1 > \beta_2$  then there are no nontrivial translation invariant stationary distributions.*

Here translation invariant means that the distribution is invariant under spatial shifts. This result and the others in this section are from Claudia Neuhauser's (1990) thesis. We conjecture that Theorem 2A holds without the assumption of translation invariance but that assumption is often difficult to remove. Note that Harris proved Theorem 1B for translation invariant initial distributions in 1974 but the general case was settled 15 years later.

Restricting our attention now to the special case  $\beta_1 = \beta_2 > \beta_c(\mathbb{Z}^d)$ , we have

**Theorem 2B.** *In dimensions  $d \leq 2$ , for any initial configuration, we have  $P(\xi_t(x) = 1, \xi_t(y) = 2) \rightarrow 0$  for all  $x, y \in \mathbb{Z}^d$ , so all stationary distributions are trivial.*

**Theorem 2C.** In dimensions  $d \geq 3$ , there is a one parameter family of stationary distributions  $v_\theta$ ,  $\theta \in [0, 1]$ , and all translation invariant stationary distributions are convex combinations of the  $v_\theta$ .

As in the voter model, (see Liggett (1985) Chapter V or Durrett (1988) Chapter 2), the dichotomy between the behavior in  $d \leq 2$  and  $d \geq 3$  comes from the fact that random walks are recurrent in the first case and transient in the second. The stationary distributions are constructed by starting the system from an initial product measure in which 1's have density  $\theta$  and 2's have density  $1 - \theta$ , i.e.,  $\xi_0(x)$  are independent and take values 1 and 2 with probabilities  $\theta$  and  $1 - \theta$ . The reader should note that while the basic contact process has a single nontrivial stationary distribution, the two color version has a one parameter family in  $d \geq 3$ .

### 3. Successional Dynamics

In this model we again have  $\xi_t : \mathbb{Z}^d \rightarrow \{0, 1, 2\}$  but this time we think of 0 as vacant and 1 and 2 as occupied by a bush or tree respectively. With this interpretation in mind the dynamics are formulated as follows:

- (i) Particles of type  $i$  die at rate one, give birth at rate  $\beta_i$ .
- (ii) A particle born at  $x$  is sent to a  $y$  chosen at random from  $\{y : \|x - y\|_1 \leq M\}$ , where  $M$  is an integer.
- (iii) If  $\xi_t(y) \geq \xi_t(x)$  then the birth is suppressed.

In words, trees can give birth onto sites occupied by bushes but not conversely. In biological terms the two species are part of a successional sequence. When only one type of particle is present, the system reduces, as in the last example, to a contact process so if  $\beta_1, \beta_2 > \beta_c$  then there are three trivial equilibria:  $\delta_0$ ,  $\mu_1$ , and  $\mu_2$ , where  $\mu_i$  is the limit starting from  $\xi_t(x) \equiv i$ .

Again, the main question to be answered is: “Are there nontrivial stationary distributions?” or more briefly “Is coexistence possible?” Our first answer is

**Theorem 3A.** If  $d = 1$  and  $M = 1$  then for any initial configuration we have  $P(\xi_t(x) = 1, \xi_t(y) = 2) \rightarrow 0$  as  $t \rightarrow \infty$  for all  $x, y \in \mathbb{Z}$  so there is no coexistence.

This result can be proved by drawing a picture of a “typical” realization of the process starting with a single 2

0010102022020002101001

and checking that since  $M = 1$  there can never be a 1 between the leftmost and rightmost 2's. If the 2's do not die out, then the ends of the interval of 2's go to  $-\infty$  and  $\infty$  respectively (see Durrett (1980)) and the 1's get crowded out. In general either (a) all the 2's die out, or (b) some 2 starts an interval that grows forever. In either case  $P(\xi_t(x) = 1, \xi_t(y) = 2) \rightarrow 0$  as  $t \rightarrow \infty$ .

We believe that coexistence is possible in all other cases

**Conjecture 3A.** *If  $d > 1$  or  $M > 1$  then coexistence is possible when  $\beta_2 = \beta_c + \varepsilon$  and  $\beta_1$  is large.*

The main trouble with proving this conjecture is that coexistence can only occur near the critical value. It is not hard to show that if  $\beta_2 > \beta(d, M)$  then there is no coexistence for any  $\beta_1 \leq \infty$ . Somewhat surprisingly, this problem which is difficult to solve when  $d = 1$  and  $M = 2$ , or  $d = 2$  and  $M = 1$  turns out to be more tractable when  $M$  is large. In addition to proving Theorem 3A, Durrett and Swindle (1990) have shown

**Theorem 3B.** *If  $\beta_1 > \beta_2^2 > 1$  then coexistence occurs for large  $M$ .*

To explain the last conclusion we need to introduce the long range contact process, a modification of the basic contact process in which (ii) is changed to:

(ii) A particle born at  $x$  is sent to a  $y$  chosen at random from  $\{y : \|x - y\|_1 \leq M\}$ .

If we write  $\beta_c(M)$  to indicate the dependence of the critical value on  $M$  and use  $\xi_\infty^1$  to denote the limit starting from all 1's then we have

**Theorem 3C.** *As  $M \rightarrow \infty$ ,  $\beta_c(M) \rightarrow 1$ . Furthermore, if  $\beta > 1$  then  $\xi_\infty^1$  converges weakly to a product measure with density  $(\beta - 1)/\beta$ .*

This result (for the neighborhood  $\{y : \|x - y\|_\infty \leq M\}$ ) was proved by Bramson, Durrett, and Swindle (1989) who identified the rate at which  $\beta_c(M)$  approached 1. A simpler and more general proof, which does not give the right rate, can be found in Durrett (1989).

To explain the condition in Theorem 3B, observe that  $\eta_t = \{x : \xi_t(x) = 2\}$  is a long range contact process, so if  $M$  is large and we are in equilibrium,  $\eta_t$  is approximately a product measure with density  $(\beta_2 - 1)/\beta_2$ . If the 2's were exactly that product measure, a 1 would die at rate  $1 + \frac{\beta_2 - 1}{\beta_2} \beta_2$ , (the second term representing births onto the site by 2's) and give birth at rate  $\beta_1/\beta_2$  (the site must not be occupied by a 2 for a successful birth to occur). So for coexistence to occur we need  $1 + \frac{\beta_2 - 1}{\beta_2} \beta_2 < \beta_1/\beta_2$  or  $\beta_1 > \beta_2^2$ . The careful reader will have noted that we have just argued the condition is necessary while Theorem 3B proves it is sufficient. Having faith in the heuristic argument we make

**Conjecture 3B.** *If  $\beta_1 < \beta_2^2$  then there is no coexistence for large  $M$ .*

*Remark.* The heuristic argument generalizes easily to show that if the two particles die at different rates then we need

$$\delta_1 + \frac{\beta_2 - \delta_2}{\beta_2} \beta_2 > \beta_1 \frac{\delta_2}{\beta_2}$$

and the proof of Theorem 3B generalizes to show that this condition is sufficient. It is natural to generalize the multitype contact process in this way but we do not know how to prove any results in that generality. The naive guess is that

$\beta_1/\delta_1 > \beta_2/\delta_2$  is right hypothesis for Theorem 2A. We believe this is correct but have no idea how to prove it.

Having discussed the existence of nontrivial stationary distributions, we turn to the question of uniqueness. Durrett and Møller (1991) have proved a “complete convergence theorem.” To state their result let  $\delta_0$ ,  $\mu_1$ , and  $\mu_2$  be the trivial stationary distributions mentioned at the beginning of this section. Let  $\mu_{12}$  be the nontrivial stationary distribution constructed in Theorem 3B. Let  $\eta_t = \{x : \xi_t(x) = 1\}$ ,  $\zeta_t = \{x : \xi_t(x) = 2\}$ ,  $\tau_1 = \inf\{t : \eta_t = \emptyset\}$ , and  $\tau_2 = \inf\{t : \zeta_t = \emptyset\}$ .

**Theorem 3C.** *If  $\beta_1 > \beta_2^2 > 1$  and  $M$  is large then*

$$\begin{aligned} \xi_t \Rightarrow P(\tau_1 < \infty, \tau_2 < \infty) \delta_0 + P(\tau_1 = \infty, \tau_2 < \infty) \mu_1 \\ + P(\tau_1 < \infty, \tau_2 = \infty) \mu_2 + P(\tau_1 = \infty, \tau_2 = \infty) \mu_{12}. \end{aligned}$$

In words if the 1’s and/or 2’s die out we end up with a trivial stationary distribution in which one or zero types of particles are present. If both the 1’s and 2’s survive and  $t$  is large, the system looks like  $\mu_{12}$  so that is the only nontrivial stationary distribution. The value of  $M$  required for Theorem 3C is larger than that for Theorem 3B which is enormous. With more work this difference might be eliminated but the interesting problem is to show

**Conjecture 3C.** *The complete convergence theorem holds whenever coexistence occurs.*

## 4. An Epidemic Model

Our fourth system is a process  $\xi_t : \mathbb{Z}^2 \rightarrow \{0, 1, 2\}$  that has been used to model the spread of epidemics and forest fires. In the epidemic interpretation 0 = healthy, 1 = infected, 2 = removed = immune or dead. In the forest fire interpretation, 0 = alive, 1 = on fire, and 2 = burnt. With these interpretations in mind, we formulate the dynamics as follows:

- (i) A burning tree sends out sparks at rate  $\beta$ .
- (ii) A spark emitted from  $x$  flies to one of the four nearest neighbors  $\{y : \|y - x\|_1 = 1\}$  chosen at random. If the spark hits a live tree, the tree catches fire and begins immediately to emit sparks.
- (iii) A tree remains on fire for an exponential amount of time with mean 1 then becomes burnt.
- (iv) Burnt trees come back to life at rate  $\alpha$ .

At first glance, the spontaneous re-appearance of trees may not seem reasonable. In the epidemic interpretation this is quite natural, however. Consider a

disease like measles that upon recovery confers lifetime immunity. New susceptibles are born and immune individuals die. We combine the two transitions into the one in (iv) to keep a constant population size.

When  $\alpha = \infty$ , sites change instantaneously from 2 to 0 and the result is the contact process. At the other extreme,  $\alpha = 0$ , is the so-called “spatial epidemic with removal” in which regrowth is impossible. We begin by considering the behavior of our processes starting with a single burning tree at the origin in the midst of an otherwise virgin forest, i.e.,  $\xi_0^0(0) = 1$ ,  $\xi_0^0(x) = 0$  for  $x \neq 0$ . Let  $\eta_t^0 = \{x : \xi_t^0(x) = 1\}$ , let  $\zeta_t^0 = \{x : \xi_t^0(x) = 2\}$ , and define a critical value by

$$\beta_c(\alpha) = \inf\{\beta : P(\zeta_t^0 \neq \emptyset \text{ for all } t) > 0\}.$$

Cox and Durrett (1988) considered the case  $\alpha = 0$  and showed

**Theorem 4A.** *If  $\beta > \beta_c(0)$  then there is a nonrandom convex set  $D$  so that on  $\{\eta_t \neq \emptyset \text{ for all } t\}$  we have  $\zeta_t^0 \approx \zeta_\infty^0 \cap tG$ , and  $\eta_t^0 \approx t\partial G$ . To be precise, for any  $\varepsilon > 0$  the following inequalities hold for large  $t$*

$$\begin{aligned}\zeta_\infty^0 \cap (1 - \varepsilon)tG &\subset \zeta_t^0 \subset (1 + \varepsilon)tG \\ \eta_t^0 &\subset (1 + \varepsilon)tG - (1 - \varepsilon)tG.\end{aligned}$$

In words, this result says that the fire expands linearly and has an asymptotic shape. The statement is made contorted by the fact that the set of trees that will ever burn,  $\zeta_\infty^0$ , is not all of  $\mathbb{Z}^d$ . Thus what we prove is that when  $t$  is large,  $\zeta_t^0$  is contained in  $(1 + \varepsilon)tG$  and (if nonempty) contains all the points of  $\zeta_\infty^0$  in  $(1 - \varepsilon)tG$ .

When  $\alpha = 0$  the system cannot have a nontrivial stationary distribution but Durrett and Neuhauser (1991) have shown

**Theorem 4B.** *If  $\beta > \beta_c(0)$  and  $\alpha > 0$  then there is a nontrivial stationary distribution, i.e., one that assigns no mass to “all healthy” state.*

The last result illustrates some of the frustrations in “applied probability.” The proof is intricate and required several months to put down on paper, but we have been repeatedly told by physicists and biologists that the conclusion is obvious. In view of our difficulties in proving existence the reader should not be surprised to learn that we have little to say about uniqueness.

**Conjecture 4C.** *If  $\beta > \beta_c(\alpha)$  then there is a unique nontrivial stationary distribution.*

In the first three examples we have had varying degrees of success in identifying the set of stationary distributions. In each of those cases however there is a useful “duality equation” and we have not been able to find one here.

## References

- Bezuidenhout, C., Grimmett, G. (1990): The critical contact process dies out. *Ann. Probab.* (to appear)
- Bramson, M., Durrett, R., Swindle, G. (1989): Statistical mechanics of crabgrass. *Ann. Probab.* **17**, 444–481
- Brower, R.C., Furman, M.A., Moshe, M. (1978): Critical exponents for the Reggeon quantum spin model. *Phys. Letters* **76B**, 213–219
- Cox, J.T., Durrett, R. (1988): Limit theorems for the spread of epidemics and forest fires. *Stoch. Processes Appl.* **30**, 171–191
- Durrett, R. (1980): On the growth of one-dimensional contact processes. *Ann. Probab.* **8**, 890–907
- Durrett, R. (1988): Lecture Notes on Particle Systems and Percolation. Wadsworth Pub. Co., Pacific Grove, CA
- Durrett, R. (1989): A new method for proving the existence of phase transitions. To appear in the Proceedings of a Conference in Honor of Ted Harris. Birkhäuser, Boston
- Durrett, R., Griffeath, D. (1982): Contact processes in several dimensions. *Z. Wahrsch. Verw. Gebiete* **59**, 535–552
- Durrett, R., Møller, A.M. (1991): Complete convergence theorem for a competition model. *Z. Wahrsch. Verw. Gebiete* (to appear)
- Durrett, R., Neuhauser, C. (1991): Epidemics with recovery in  $d = 2$ . *Adv. Appl. Probab.* (to appear)
- Durrett, R., Schonmann, R.H. (1987): Stochastic growth models. In: Kesten (1987)
- Durrett, R., Swindle, G. (1990): Are there bushes in a forest? *Stoch. Processes Appl.* (to appear)
- Griffeath, D. (1978): Limit theorems for non-ergodic set-valued Markov processes. *Ann. Probab.* **6**, 379–387
- Griffeath, D. (1983): The binary contact path process. *Ann. Probab.* **11**, 692–705
- Harris, T.E. (1974): Contact interactions on a lattice. *Ann. Probab.* **2**, 969–988
- Holley, R., Liggett, T.M. (1978): The survival of contact processes. *Ann. Probab.* **6**, 198–206
- Holley, R., Liggett, T.M. (1981): Generalized potlatch and smoothing processes. *Z. Wahrsch. Verw. Gebiete* **55**, 165–195
- Kesten, H. (ed.) (1987): Percolation theory and the ergodic theory of interacting particle systems. Springer, New York Berlin Heidelberg
- Liggett, T.M. (1985): Interacting particle systems. Springer, New York Berlin Heidelberg
- Neuhauser, C. (1990): Ergodic theorems for the multitype contact process. Ph.D. Thesis, Cornell University

# Recurrent Ergodic Structures and Ramsey Theory

Hillel Furstenberg

Institute of Mathematics, Hebrew University, Givat Ram, Jerusalem, Israel

## Introduction

Ramsey theory is best defined by example and the classic example of a Ramsey type theorem is the result of van der Waerden: if the integers  $\mathbb{Z}$  are partitioned into finitely many sets, one of these contains arbitrarily long arithmetic progressions. An equivalent version states: given natural numbers  $r, l$ ,  $\exists N(r, l)$  so that if  $N \geq N(r, l)$  and the integers  $\{1, 2, \dots, N\}$  are partitioned into  $r$  sets, one of these contains an  $l$ -term arithmetic progression. Erdős and Turán realized that this result would follow if it were the case that *any* subset of  $\{1, 2, \dots, N\}$  comprising  $\delta N$  elements contains an  $l$ -term arithmetic progression provided  $N$  is sufficiently large. This result, conjectured in 1936 [ET1] was proved by E. Szemerédi in 1973 [Sz1], and was reproved using ergodic theoretic methods in 1976 [Fu1]. The theorems of van den Waerden and Szemerédi illustrate the two sides of Ramsey theory: coloring – or partition – results, and density results. The latter are clearly stronger than the former, and the proofs are typically more recondite. The role of ergodic theory in density theorems of this type stems from the fact that in a number of situations theorems about patterns found in sets having “density” bounded from below are equivalent to theorems about the “return”, or “recurrence”, patterns for measure preserving transformations acting on a measure space. It would appear that the ubiquity of certain patterns in the combinatorics of large sets reflects the phenomenon of recurrence in ergodic theoretic contexts and the latter has to be studied to gain insight into the former.

We shall be examining three different contexts for recurrence results in ergodic theory. We shall mention the combinatorial (Ramsey theoretic) equivalents of these results, although we shall have to refer the reader elsewhere for the proof of the equivalence. Our purpose here is to display the common features of these recurrence phenomena. We shall find that in each setting there is a notion of *rigidity* and the notion of a special system constructed from scratch by a finite succession of rigid extension. For these systems (*distal* and *quasi-distal*) it will be possible to deduce recurrence properties directly, the key tool being Ramsey theory of the van der Waerden sort. For arbitrary systems we shall find a method to “factor out” the “independent” component, reducing the recurrence property

for a general system to that of a factor system, which will be one of the special systems for which the result has been established.

To motivate the ergodic theoretic discussion we list here some of the combinatorial consequences. For more details, particularly as regards the connection between the ergodic theoretic formulation and the Ramsey theoretic results, we refer the reader to [Fu2, FK2, FK4, and GRS1].

**Theorem A.** *There is a function  $N(\delta, l)$  for  $\delta > 0$  and  $l \in \mathbb{N}$  so that if  $N \geq N(\delta, l)$  and  $S \subset \{1, 2, \dots, N\}$  with  $|S| \geq \delta N$  then  $S$  contains an  $l$ -term arithmetic progression. ( $|S|$  denotes cardinality of  $S$ .)*

**Theorem B.** *Let  $P$  denote a finite subset of  $\mathbb{R}^d$ . There is a function  $R(\delta, P)$  so that if  $R \geq R(\delta, P)$  and  $B$  is a ball of radius  $R$  in  $\mathbb{R}^d$  and  $S \subset B$  is a measurable subset with  $m(S) \geq \delta m(B)$ , then  $S$  contains an integral dilation of  $P$ , i.e., a set of the form  $u + nP$ ,  $u \in \mathbb{R}^d$ ,  $n \in \mathbb{N}$ .*

**Theorem C.** *There exists a function  $D(\delta, d, q)$  for  $\delta > 0$ ,  $d \in \mathbb{N}$ ,  $q = a$  prime power, so that if  $V$  is a vector space over  $\mathbb{E}_q$  of dimension  $\geq D(\delta, d, q)$  and  $S \subset V$  is a subset with  $|S| \geq \delta |V|$ , then  $S$  contains a  $d$ -flat (translate of  $d$ -dimensional subspace).*

**Theorem D.** *If  $A$  is a finite set (alphabet) let  $W_N(A)$  denote words of length  $N$  in  $A$ . Let  $W_N^*(A) = W_N(A \cup \{t\}) \setminus W_N(A)$  so that  $\varphi(t) \in W_N^*(A)$  is a word in which the variable  $t$  appears. There exists a function  $L(\delta, |A|)$  so that if  $N \geq L(\delta, |A|)$  and  $S \subset W_N(A)$  with  $|S| \geq \delta |W_N(A)|$ , then there exists  $\varphi(t) \in W_N^*(A)$  so that the words  $\varphi(a), \varphi(b), \dots, \varphi(s)$ ,  $a, b, \dots, s \in A$  are all in  $S$ .*

Theorem A is Szemerédi's theorem. Theorem B is equivalent to a multidimensional version of Szemerédi's theorem. It was proved first in [FK1]. Theorem C is a density theoretic version of a coloring theorem first conjectured by Roth and proved by R.L. Graham and B. Rothschild [GR1]. The latter also follows from the theorem of A.W. Hales and R.I. Jewett [HJ1]. A special case of Theorem C was proved by T.C. Brown and J.P. Buhler [BB1], and in generality it was proved in [FK2]. Theorem D was conjectured by R.L. Graham and proved in [FK4].

By interpreting the letters of  $A$  as the digits for writing numbers to a base  $b$ , we can see that  $D \Rightarrow A$ . On the other hand, interpreting  $A$  as a vector space over a finite field, we can see that  $D \Rightarrow C$ . It is also not difficult to show that  $D \Rightarrow B$ , so that Theorem D holds the key to a number of results of Ramsey theory.

We remark that the major part of what is presented here is joint work with Y. Katznelson.

## 1. Classical Systems

In the classical setting for ergodic theory we speak of a *measure preserving system* (m.p.s.)  $(X, \mathcal{B}, \mu, T)$  where  $T : X \rightarrow X$  is a measurable, measure preserving map, and  $\mu$  is a probability measure on the sets of  $\mathcal{B}$ . It is useful to distinguish between

*ergodic* and *non-ergodic* systems.  $(X, \mathcal{B}, \mu, T)$  is *ergodic* if  $T^{-1}A = A$  implies  $\mu(A) = 0$  or 1. There is a rather general decomposition theorem for measure preserving systems to ergodic components, and the recurrence phenomena we shall study here follow for arbitrary systems once they are established for ergodic systems. We also can assume without loss of generality that  $T$  is invertible.

If  $T$  is a measure preserving transformation on  $X$  then  $T$  induces an operator on functions on  $X$ :  $T^{-1}f(x) = f(Tx)$ . (The notation corresponds to the fact that  $S^{-1}(T^{-1}f) = (TS)^{-1}f$ .)  $T^{-1}$  defines a unitary operator on  $L^2(X, \mathcal{B}, \mu)$  and we can study spectral properties of  $T^{-1}$ . In particular ergodicity is equivalent to 1 being a simple eigenvalue, and if, in addition, there are no other eigenvalues, we say that  $T$ , or  $(X, \mathcal{B}, \mu, T)$ , is *weak mixing*.

The mean ergodic theorem for  $L^2$ -functions is a special case of a general theorem for contraction operators in Hilbert space:

**Theorem 1.** *If  $(X, \mathcal{B}, \mu, T)$  is an ergodic system and  $f \in L^2(X, \mathcal{B}, \mu)$ , then*

$$\frac{f + T^{-1}f + T^{-2}f + \cdots + T^{-N}f}{N+1} \rightarrow \int f d\mu$$

*in the norm of  $L^2(X, \mathcal{B}, \mu)$ .*

V. Bergelson has given a far reaching generalization of this in the case of a weak mixing system [Be1].

**Theorem 2.** *Let  $(X, \mathcal{B}, \mu, T)$  be a weak mixing system. Let  $p_1(t), p_2(t), \dots, p_k(t)$  be polynomials with integer coefficients with  $\deg p_i > 0$  and  $\deg(p_i - p_j) > 0$  for  $i \neq j$ . Then if  $f_1, f_2, \dots, f_k \in L^\infty(X, \mathcal{B}, \mu)$  we will have*

$$\frac{1}{N+1} \sum_{n=0}^N T^{-p_1(n)} f_1 T^{-p_2(n)} f_2 \cdots T^{-p_k(n)} f_k \rightarrow \int f_1 d\mu \int f_2 d\mu \cdots \int f_k d\mu \quad (1)$$

*in the norm of  $L^2(X, \mathcal{B}, \mu)$ .*

One consequence of this theorem is that if  $(X, \mathcal{B}, \mu, T)$  is weak mixing and  $A_1, A_2, \dots, A_k$  are any  $k$  subsets of positive measure, we will have for some (in fact, most)  $n$ :

$$\mu(A_1 \cap T^{-n}A_2 \cap \cdots \cap T^{-(k-1)n}A_k) > 0, \quad (2)$$

and, in particular, if  $\mu(A) > 0$ ,

$$\mu(A \cap T^{-n}A \cap \cdots \cap T^{-(k-1)n}A) > 0. \quad (3)$$

The foregoing property, *that points of  $A$  return to  $A$  along arithmetic progressions*, is the recurrence phenomenon we have in mind in the case of a classical system. In the weak mixing case we have found that it follows from a generalized ergodic theorem.

The proof of Theorem 2 is based on the following lemma for convergence of averages in Hilbert space to 0.

**Lemma 3.** Let  $\{u_n\}$  be a bounded sequence of vectors in a Hilbert space. Assume for each  $m$ , that

$$\gamma_m = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \langle u_n, u_{n+m} \rangle$$

exists, and that

$$\frac{1}{M} \sum_{m=1}^M \gamma_m \rightarrow 0 .$$

then

$$\left\| \frac{1}{N} \sum_{n=1}^N u_n \right\| \rightarrow 0 .$$

Using the lemma and the ergodic theorem, one proves (1.1) for progressively more complex situations, beginning with

$$\frac{1}{N+1} \sum_{n=0}^N T^{-an} f T^{-bn} g, \quad \frac{1}{N+1} \sum_{n=0}^N T^{-an} f T^{-bn} g T^{-cn} h, \dots$$

We refer to [Be1] for details.

When  $(X, \mathcal{B}, \mu, T)$  is not weak mixing the result of Theorem 2 cannot be expected to be valid. For example assume  $\varphi$  is an eigenfunction,  $T^{-1}\varphi = \lambda\varphi$ , and set  $f = \varphi^2, g = \varphi^{-1}$ , then  $T^{-n}f T^{-2n}g = (\lambda^{2n}\varphi^2)(\lambda^{-2n}\varphi^{-1}) = \varphi$  and

$$\frac{1}{N+1} \sum_0^N T^{-n} f T^{-2n} g \rightarrow \varphi$$

which is not constant. In the weak mixing case we would have obtained  $\int f d\mu \int g d\mu$  as the limit.

## 2. Isometric Extensions and Distal Factors

If  $(X, \mathcal{B}, \mu, T)$  and  $(Y, \mathcal{D}, \nu, T)$  are two systems (it is convenient to use the same letter for the transformation), we say that a map  $\pi : X \rightarrow Y$  defines  $(Y, \mathcal{D}, \nu, T)$  as a factor of  $(X, \mathcal{B}, \mu, T)$  if  $\pi$  is measurable, measure preserving, and  $\pi(Tx) = T(\pi x)$ . We lift functions from  $Y$  to  $X$ ,  $f \rightarrow f \circ \pi$  and since this imbeds  $L^2(Y, \mathcal{D}, \nu)$  isometrically into  $L^2(X, \mathcal{B}, \mu)$ , we shall identify  $f \in L^2(Y, \mathcal{D}, \nu)$  with  $f \circ \pi \in L^2(X, \mathcal{B}, \mu)$ . It will be convenient to denote the conditional expectation operator  $f \mapsto E(f | \pi^{-1}\mathcal{D})$  by  $f \mapsto E(f | Y)$  and we regard  $E(f | Y)$  either as a function on  $Y$  or on  $X$ . Note that  $T^{-1}$  commutes with  $E(\cdot | Y)$ .

We call a m.p.s. Kronecker if it has the form  $(Z, \mathcal{D}, m, T)$  where  $(Z, +)$  is a compact abelian group with a dense cyclic subgroup  $\mathbb{Z}\alpha$ ,  $\mathcal{D} =$  Borel sets,  $m =$  Haar measure, and  $Tz = z + \alpha$ . For a Kronecker system,  $L^2(Z, \mathcal{D}, m)$  is spanned by eigenfunctions (each character is an eigenfunction). It follows that if a system  $(X, \mathcal{B}, \mu, T)$  possesses a non-trivial Kronecker factor then it cannot be

weak mixing. The converse is also true: if  $(X, \mathcal{B}, \mu, T)$  is not weak mixing it has a non-trivial Kronecker factor.

We describe a relativized notion of a Kronecker factor, representing a certain relationship between a system and a factor.

**Definition 1.**  $(X, \mathcal{B}, \mu, T)$  is an *isometric extension* of a factor  $(Y, \mathcal{D}, \nu, T)$  if it can be represented – up to isomorphism – as follows:  $X = Y \times M$ ,  $M$  is a homogeneous space of a compact group  $G$ ,  $\mu = \nu \times m$  where  $m$  is the  $G$ -invariant measure on  $M$ , and there exists a measurable function  $\varrho : Y \rightarrow G$  so that for  $y \in Y, u \in M$ ,

$$T(y, u) = (Ty, \varrho(y)u).$$

An ergodic isometric extension of the trivial one-point system is a Kronecker system.

**Definition 2.** A m.p.s.  $(X, \mathcal{B}, \mu, T)$  is *(n-step) distal* if it has a succession of factors

$$(X, \mathcal{B}, \mu, T) = (X_1, \mathcal{B}_1, \mu_1, T) \rightarrow \cdots \rightarrow (X_n, \mathcal{B}_n, \mu_n, T) \rightarrow \text{1-point} \quad (4)$$

where each extension is isometric.

The terminology comes from topological dynamics where distal systems are defined by the property that for  $x \neq x'$ , the distance  $d(T^n x, T^n x')$  is bounded from below. It is easily seen that for metric spaces, a succession of isometric extensions leads to a distal system in this sense.

We now formulate the generalization of V. Bergelson's theorem 1.2 to the general ergodic situation. Note that weak mixing systems have only trivial distal factors.

**Theorem 3.** Let  $(X, \mathcal{B}, \mu, T)$  be an ergodic m.p.s. Let  $p_1(t), p_2(t), \dots, p_k(t)$  be polynomials with integer coefficients with  $\deg p_i > 0$  and  $\deg(p_i - p_j) > 0$  for  $i \neq j$ . There exists a distal factor  $\pi : X \rightarrow Y$  so that for any  $f_1, f_2, \dots, f_k \in L^\infty(X, \mathcal{B}, \mu)$

$$\begin{aligned} & \frac{1}{N+1} \sum_{n=0}^N T^{-p_1(n)} f_1 T^{-p_2(n)} f_2 \cdots T^{-p_k(n)} f_k \\ & - \frac{1}{N+1} \sum_{n=0}^N T^{-p_1(n)} E(f_1 \mid Y) \cdots T^{-p_k(n)} E(f_k \mid Y) \rightarrow 0 \end{aligned} \quad (5)$$

in  $L^2(X, \mathcal{B}, \mu)$ .

We call  $(Y, \mathcal{D}, \nu, T)$  a *characteristic factor* for  $(T^{-p_1(n)}, \dots, T^{-p_k(n)})$ . We use the indefinite article here because any extension of a characteristic factor is again a characteristic factor. We do not know if there exists a “smallest” one.

For special cases the information is more specific [Fu3]:

**Theorem 4.** For any ergodic m.p.s.  $(X, \mathcal{B}, \mu, T)$  there exists a “largest” Kronecker factor  $(Z, \mathcal{D}, m, T)$  and for  $f, g \in L^\infty(X, \mathcal{B}, \mu)$ ,  $a, b \in \mathbb{Z}$ ,  $a \neq b$ ,

$$\lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N T^{-an} f T^{-bn} g = \lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N T^{-an} E(f | Z) T^{-bn} E(g | Z) . \quad (6)$$

Here we are also guaranteed existence of the limits in question (in  $L^2(X, \mathcal{B}, \mu)$ ). If we write  $\bar{f} = E(f | Z)$ ,  $\bar{g} = E(g | Z)$  and regard the right hand side of (2.3) as a function  $\psi(\zeta)$  on  $Z$ , then it can be rewritten as

$$\psi(\zeta) = \int \bar{f}(\zeta + a\theta) \bar{g}(\zeta + b\theta) d\theta$$

We can also be more explicit regarding characteristic factors for  $(T^{-an}, T^{-bn}, T^{-cn})$ . To formulate this result we need the notion of a *2-step nilpotent system*. This is a system for which the underlying space has the form  $N/\Gamma$  where  $N$  is a 2-step nilpotent group,  $\Gamma$  is a closed cocompact subgroup and  $N$  possesses a sequence of closed subgroups  $N_i$  so that  $N/N_i$  is locally compact,  $\Gamma/\Gamma \cap N_i$  is discrete and cocompact in  $N/N_i$ , the measure on  $N/\Gamma$  is the (unique)  $N$ -invariant probability measure (lifted from  $(N/N_i) / (\Gamma/\Gamma \cap N_i)$ ) and the measure preserving transformation is  $T(g\Gamma) = (\tau g)\Gamma$  for some  $\tau \in N$ .

**Theorem 5.** For any ergodic  $(X, \mathcal{B}, \mu, T)$  there exists a 2-step nilpotent factor  $(Y, \mathcal{D}, v, T)$  so that for  $f, g, h \in L^\infty(X, \mathcal{B}, \mu)$  and  $a, b, c$  distinct integers, we have

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N T^{-an} f T^{-bn} g T^{-cn} h = \\ \lim_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N T^{-an} E(f | Y) T^{-bn} E(g | Y) T^{-cn} E(h | Y) \end{aligned}$$

Theorem 5 was obtained jointly with B. Weiss and overlaps results of Conze and Lesigne [CL1]. It may be conjectured that similar results are valid for more general expressions with 2-step nilpotent systems replaced by arbitrary nilpotent systems.

### 3. Multiple Recurrence

We begin this section by formulating a version of (1.3) for distal systems.

**Theorem 1.** Let  $(Y, \mathcal{D}, v, T)$  be a distal system, let  $f \in L^\infty(Y, \mathcal{D}, v)$  be non-negative but not vanishing a.e. Then for any  $l \in \mathbb{N}$ ,

$$\liminf_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N \int f T^{-n} f T^{-2n} f \cdots T^{-(l-1)n} f d\mu > 0 . \quad (7)$$

To explain the role of distality in this result we shall prove here by way of illustration that for 1-step distal systems one has the weaker statement: With  $(Y, \mathcal{D}, v, T)$  and  $f, l$  as above,  $\exists n$  such that

$$\int f T^{-n} f T^{-2n} f \cdots T^{-(l-1)n} f d\mu > 0 . \quad (8)$$

A refinement of the argument will show that (7) is valid. Moreover a relativized version of this argument can be reapplied successively to handle the general distal case.

If  $M$  is a homogeneous space of the compact group  $G$ , then  $G$  acts continuously on  $L^2(M, m) : \psi(u) \rightarrow \psi^g(u) = \psi(gu)$ . In particular, if  $\psi \in L^\infty(M, m)$ , its orbit under  $G$ ,  $\Psi = \{\psi^g \mid g \in G\}$ , is compact in  $L^2(M, m)$ . Suppose  $\psi \geq 0$  and that  $\psi$  is not 0 a.e. For any  $l$ ,  $\exists \varepsilon > 0$  so that if  $\psi_1, \psi_2, \dots, \psi_l \in L^\infty(M, m)$  satisfy  $\|\psi_j - \psi\| < \varepsilon$  then  $\int \psi_1 \psi_2 \cdots \psi_l dm > 0$ . Since  $\Psi$  is compact we can partition  $\Psi$  into finitely many subsets of diameter  $< \varepsilon$ . So we can write  $\Psi = \Psi_1 \cup \cdots \cup \Psi_r$  so that if  $\psi_1, \psi_2, \dots, \psi_l \in \Psi_j$  then  $\int \psi_1 \psi_2 \cdots \psi_l dm > 0$ .

Now consider a 1-step distal system, i.e., a Kronecker system,  $(Z, \mathcal{D}, v, T)$ .  $Z$  plays the role of  $G$  as well as  $M$ . Take  $f$  as in the theorem and set  $\psi = f$ . By van der Waerden's theorem there is an  $l$ -term arithmetic progression,  $a, a+n, a+2n, \dots, a+(l-1)n$  such that each  $T^{-(a+in)} f, i = 0, 1, \dots, l-1$ , belongs to the same  $\Psi_j$ . Then

$$\int T^{-a} f T^{-(a+n)} f \cdots T^{-(a+(l-1)n)} f d\mu > 0$$

or

$$\int f T^{-n} f \cdots T^{-(l-1)n} f d\mu > 0$$

as claimed.

We now combine Theorem 1 with Theorem 2.1 to obtain

**Theorem 2.** *If  $(X, \mathcal{B}, \mu, T)$  is any ergodic m.p.s.,  $f \in L^\infty(X, \mathcal{B}, \mu)$  is a non-negative function not vanishing a.e., then*

$$\liminf_{N \rightarrow \infty} \frac{1}{N+1} \sum_{n=0}^N \int f T^{-n} f T^{-2n} f \cdots T^{-(l-1)n} f d\mu > 0 .$$

Decomposition to ergodic components gives the result for arbitrary measure preserving systems. Finally taking  $f = 1_A$  we obtain the

**Theorem (Multiple Recurrence).** *Let  $(X, \mathcal{B}, \mu, T)$  be any m.p.s., let  $A \in \mathcal{B}$  with  $\mu(A) > 0$ . Then for any  $l$  there exists  $n$  with*

$$\mu(A \cap T^{-n} A \cap T^{-2n} A \cap \cdots \cap T^{-(l-1)n} A) > 0 . \quad (9)$$

As we show in [Fu1] and [Fu2] this is equivalent to Theorem A of the Introduction.

## 4. IP-Systems

We denote by  $\mathcal{F}$  the collection of all finite subsets of  $\mathbb{N} = \{1, 2, 3, \dots\}$ . If  $\alpha, \beta \in \mathcal{F}$  we write  $\alpha < \beta$  if  $\max \alpha < \min \beta$ .

**Definition 1.** An *IP-subset* of  $\mathcal{F}$  consists of a sequence  $\alpha_1 < \alpha_2 < \dots$  in  $\mathcal{F}$  together with all finite unions of these.

The elements of an IP-subset of  $\mathcal{F}$  may be placed in one-one correspondence with  $\mathcal{F}: \beta = \{i_1, i_2, \dots, i_k\} \leftrightarrow \alpha_{i_1} \cup \alpha_{i_2} \cup \dots \cup \alpha_{i_k} = \alpha_\beta$ .

**Definition 2.** An  *$\mathcal{F}$ -sequence* is any sequence  $\{x_\alpha\}$  labelled by  $\alpha \in \mathcal{F}$ . An  *$\mathcal{F}$ -subsequence* of  $\{x_\alpha\}$  is the subsequence corresponding to an IP-subset of  $\mathcal{F}$ :  $x'_\beta = x_{\alpha_\beta}$ .

**Definition 3.** If  $X$  is a topological space,  $\{x_\alpha\}$  an  $\mathcal{F}$ -sequence in  $X$ , we say that  $x_\alpha \rightarrow x_0 \in X$ , or  $\lim_{\alpha \rightarrow \infty} x_\alpha = x_0$ , if for any neighborhood  $V$  of  $x_0$ ,  $\exists \alpha_0$  so that for  $\alpha > \alpha_0$  we have  $x_\alpha \in V$ .

The following theorem is equivalent to N. Hindman's theorem on finite sum sets [Hi1]:

**Theorem 4.** If  $X$  is a compact metric space and  $\{x_\alpha\}$  is an  $\mathcal{F}$ -sequence in  $X$ , then  $\{x_\alpha\}$  has a convergent  $\mathcal{F}$ -subsequence.

If  $X$  is finite this can be reformulated as: if  $\mathcal{F} = \mathcal{C}_1 \cup \dots \cup \mathcal{C}_r$  then some  $\mathcal{C}_j$  contains an IP-subset of  $\mathcal{F}$ . This latter formulation is due to Hindman.

The next definition describes the principal concept of this section.

**Definition 5.** An  $\mathcal{F}$ -sequence  $\{\sigma_\alpha\}$  with values in a semigroup is an *IP-system* if  $\alpha < \beta$  implies that

$$\sigma_{\alpha \cup \beta} = \sigma_\alpha \sigma_\beta . \quad (10)$$

Writing  $\sigma_i$  for  $\sigma_{\{i\}}$  we see that  $\sigma_{\{i_1, i_2, \dots, i_k\}} = \sigma_{i_1} \sigma_{i_2} \dots \sigma_{i_k}$  when  $i_1 < i_2 < \dots < i_k$ . We note that in the commutative case (4.1) is valid whenever  $\alpha \cap \beta = \emptyset$ .

**Lemma 6.** An  $\mathcal{F}$ -subsequence of an IP-system is an IP-system.

We call this an *IP-subsystem*.

The classical notion of a m.p.s. will now be replaced by that of an IP-system of measure preserving transformations (m.p.t.). Inside a classical system we can form numerous IP-systems, taking  $T_\alpha = T^{n_\alpha}$  where  $n_\alpha = \sum_{i \in \alpha} n_i$  is an IP-system in  $(\mathbb{Z}, +)$ . In addition to providing more information regarding classical systems, the analysis of IP-systems provides recurrence results for situations where the classical setup doesn't prevail – or is trivial, because as might be the case  $T^p$  could equal the identity for some  $p$  and all  $T$  under consideration.

The main theorem of this section is

**Theorem (IP-Recurrence).** Let  $(X, \mathcal{B}, \mu)$  be a measure space and let  $\{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\}$  be  $k$  IP-systems of m.p.t. of  $(X, \mathcal{B}, \mu)$  where all the  $T_\alpha^i$  commute with one another. Then if  $A \in \mathcal{B}, \mu(A) > 0, \exists \alpha \in \mathcal{F}$  so that

$$\mu(T_\alpha^{(1)-1}A \cap T_\alpha^{(2)-1}A \cap \dots \cap T_\alpha^{(k)-1}A) > 0. \quad (11)$$

We remark that this theorem implies Theorems B and C of the Introduction. We refer the reader to [FK2] for details.

The proof of the IP-recurrence theorem depends on a development of ergodic theory for IP-systems analogous to classical theory. We proceed to outline the salient points of such a theory. One important point should be made at the outset. Since in the recurrence theorem we are looking for one index  $\alpha$  for which we have a multiple intersection, we may pass from the IP-systems to IP-subsystems that may be more convenient. Since on account of Theorem 4, an  $\mathcal{F}$ -subsequence of an  $\mathcal{F}$ -sequence in a compact metric space converges, we shall often find it convenient to deal with an appropriately chosen  $\mathcal{F}$ -subsequence, which means we are really dealing with IP-subsystems of the initially given systems.

If  $\{T_\alpha\}$  is an IP-system of m.p.t. of  $(X, \mathcal{B}, \mu)$  then writing as before  $T_\alpha^{-1}f(x) = f(T_\alpha x)$  we obtain, in the commutative case, an IP-system of isometry operators  $\{T_\alpha^{-1}\}$  on  $L^2(X, \mathcal{B}, \mu)$ . Let  $\mathcal{K}$  be the semigroup of linear operators on  $L^2$  with norm  $\leq 1$ , endowed with the weak operator topology.  $\mathcal{K}$  is a compact metrizable space. Using Theorem 4 we can find a subsystem so that  $T_\alpha^{-1}$  converges in  $\mathcal{K}$ . We now have

**Proposition 7.** If  $\{K_\alpha\} \subset \mathcal{K}$  is an IP-system and  $K_\alpha \rightarrow Q$  then  $Q$  is a self-adjoint projection operator. When  $K_\alpha = T_\alpha^{-1}$  for m.p.t.  $\{T_\alpha\}$ , then  $Q$  is the conditional expectation operator relative to a subalgebra  $\mathcal{D} \subset \mathcal{B}$ .

*Proof.* If  $Q = \lim_{\alpha \rightarrow \infty} K_\alpha$  then  $Q^2 = \lim_{\alpha \rightarrow \infty} \lim_{\beta \rightarrow \infty} K_\alpha K_\beta = \lim_{\alpha \rightarrow \infty} \lim_{\beta \rightarrow \infty} K_{\alpha \cup \beta} = Q$ . The functions in the range of  $Q$  are characterized by:  $K_\alpha f \rightarrow f$  weakly. Since  $\|f\| \geq \|K_\alpha f\|$ , weak convergence implies strong convergence. For strong convergence  $K_\alpha f \rightarrow f, K_\alpha g \rightarrow g \Rightarrow K_\alpha f \vee K_\alpha g \rightarrow f \vee g$ . Now  $T_\alpha^{-1}f \vee T_\alpha^{-1}g = T_\alpha^{-1}(f \vee g)$  and so if  $f, g \in \text{range } Q$  also  $f \vee g \in \text{range } Q$ . This characterizes conditional expectation operators.  $\square$

For an IP-system  $\{T_\alpha\}$  we now suppose having passed to a subsystem for which  $\lim T_\alpha^{-1} = Q_T$  exists. For IP-systems this operator  $Q_T$  plays the role of the ergodic average. Since  $Q_T$  is self-adjoint and the  $T_\alpha^{-1}$  are isometric this convergence will generally not be in the strong topology, and certainly not pointwise. In the extreme case when  $Q_T = \text{identity}$  then the convergence is strong and we speak of  $\{T_\alpha\}$  as a *rigid* system. An example of a rigid system is one obtained from the Kronecker system  $(Z, \mathcal{B}, \mu, T)$  setting  $T_\alpha = T^{n_\alpha}$  for  $\{n_\alpha\}$  an IP-system in  $Z$ . If  $T_\alpha \rightarrow Q_T$  it is not hard to show, applying  $T_\alpha$  to eigenfunctions of  $T$ , that  $Q_T = \text{identity}$ . In IP-theory rigid systems play the role of Kronecker systems.

The opposite extreme to rigidity occurs when  $Q_T f = E(f)$ . This means that  $\int f T_\alpha^{-1} g d\mu \rightarrow \int f d\mu \int g d\mu$  so that  $f$  and  $T_\alpha^{-1}g$  are asymptotically independent. We call this case *mixing*.

Let  $(X, \mathcal{B}, \mu)$  be a measure space,  $(Y, \mathcal{D}, v)$  a factor of  $(X, \mathcal{B}, \mu)$  such that a system  $\{T_\alpha\}$  is defined both on  $X$  and on  $Y$ , that is to say, we have  $\pi : X \rightarrow Y$  with  $\pi T_\alpha = T_\alpha \pi$ . A subspace  $\mathcal{M} \subset L^2(X, \mathcal{B}, \mu)$  is called a *Y-module* if it is closed under multiplication by  $L^\infty(Y, \mathcal{D}, v)$ . A *Y-module* is of *finite rank* if it is spanned by finitely many functions over  $L^\infty(Y, \mathcal{D}, v)$ . Finally  $\mathcal{M}$  is called  $\{T_\alpha\}$  *quasi-invariant* if for each  $g \in \mathcal{M}$  and  $\varepsilon > 0$ ,  $\exists \alpha_0$  so that for  $\alpha > \alpha_0$  the  $L^2$  distance satisfies

$$\text{dist}(T_\alpha^{-1}g, \mathcal{M}) < \varepsilon .$$

**Definition 8.**  $(X, \mathcal{B}, \mu, \{T_\alpha\})$  is a *rigid extension* of  $(Y, \mathcal{D}, v, \{T_\alpha\})$  if  $L^2(X, \mathcal{B}, \mu)$  is spanned by  $\{T_\alpha\}$  quasi-invariant  $Y$ -modules of finite rank.

Now consider commuting IP-systems  $\{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\}$  of m.p.t. on  $(X, \mathcal{B}, \mu)$ . We call  $(X, \mathcal{B}, \mu, \{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\})$  a *k-fold IP-system*.

**Definition 9.**  $(X, \mathcal{B}, \mu, \{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\})$  is *quasi-distal* if there is a sequence of  $k$ -fold IP-systems  $(X_l, \mathcal{B}_l, \mu_l, \{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\})$ ,  $l = 0, 1, 2, \dots, n$ , with  $(X_0, \mathcal{B}_0, \mu_0) = (X, \mathcal{B}, \mu)$ , and with  $X_n =$  point, and with maps  $\pi_l : X_l \rightarrow X_{l+1}$  which are measure preserving and satisfy  $\pi_l T_\alpha = T_\alpha \pi_l$  and such that for each  $l$ ,  $(X_l, \mathcal{B}_l, \mu_l, \{T_\alpha^{(i)} T_\alpha^{(j)-1}\})$  is a rigid extension of  $(X_{l+1}, \mathcal{B}_{l+1}, \mu_{l+1}, \{T_\alpha^{(i)} T_\alpha^{(j)-1}\})$  for some  $i \neq j$ .

We now have

**Theorem 10.** Given a  $k$ -fold IP-system  $(X, \mathcal{B}, \mu, \{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\})$  there is a subsystem and a factor  $(Y, \mathcal{D}, v, \{T_\alpha^{(1)}\}, \dots, \{T_\alpha^{(k)}\})$  such that the latter is quasi-distal and for any  $f_1, f_2, \dots, f_k \in L^\infty(X, \mathcal{B}, \mu)$  we have

$$\lim_{\alpha} \int T_\alpha^{(1)-1} f_1 \cdots T_\alpha^{(k)-1} f_k d\mu = \lim_{\alpha} \int T_\alpha^{(1)-1} E(f_1 | Y) \cdots T_\alpha^{(k)-1} E(f_k | Y) dv .$$

The implication of Theorem 10 for us is that the IP-recurrence theorem will follow in general if we can show that for quasi-distal systems and  $f \geq 0$  with  $f$  not vanishing a.e.,

$$\lim_{\alpha \rightarrow \infty} \int T_\alpha^{(1)-1} f T_\alpha^{(2)-1} f \cdots T_\alpha^{(k)-1} f d\mu > 0 .$$

This can indeed be shown in a manner analogous to the proof of Theorem 3.1, by successively lifting the desired property via the rigid extensions that constitute the links of the quasi-distal system. This outlines the proof of the IP-recurrence theorem.

## 5. $W(\Lambda)$ -Systems

In this section we describe the ergodic structure whose recurrence properties lead to Theorem D, the density version of the Hales-Jewett theorem. The structure will be similar to that of the previous section except that now the operators are no longer assumed to commute.

We take  $\Lambda = \{1, 2, \dots, k\}$ . Suppose we are given  $k$  sequences of m.p.t.  $\{T_n^{(1)}, \{T_n^{(2)}\}, \dots, \{T_n^{(k)}\}\}$  of a measure space  $(X, \mathcal{B}, \mu)$ . Given a word  $w \in W(\Lambda)$  we can form the transformation

$$T(w) = T_1^{w(1)} T_2^{w(2)} \cdots T_l^{w(l)} . \quad (12)$$

**Definition 1.** A  $W(\Lambda)$ -system consists of a family of m.p.t. of a measure space  $(X, \mathcal{B}, \mu)$  corresponding to  $w \in W(\Lambda)$  and formed in accordance with (12). We denote the system  $(X, \mathcal{B}, \mu, \{T(w)\}_{w \in W(\Lambda)})$ .

We can now state the

**Theorem ( $W(\Lambda)$ -Recurrence).** Let  $(X, \mathcal{B}, \mu, \{T(w)\}_{w \in W(\Lambda)})$  be a  $W(\Lambda)$ -system. If  $f \in L^\infty(X, \mathcal{B}, \mu)$  is non-negative and not 0 a.e. then  $\exists \varphi(t) \in W^*(\Lambda)$  so that

$$\int T(\varphi(1))^{-1} f T(\varphi(2))^{-1} f \cdots T(\varphi(k))^{-1} f d\mu > 0 . \quad (13)$$

The proof of this theorem is patterned after the proof of the IP-recurrence theorem of §4. There is a new feature which arises on account of the non-commutativity. This feature also reflects an aspect of Ramsey theory for  $W^*(\Lambda)$  as opposed to  $\mathcal{F}$ .

We use the following notation. If  $w \in W_N(\Lambda)$  and  $\alpha \subset \{1, 2, \dots, N\}$ , we designate by  $w_\alpha$  the word of  $W^*(\Lambda)$  in which the letters at positions of  $\alpha$  are replaced by the variable  $t$ . If we write  $\varphi(t)$  for  $w_\alpha$  then we denote  $\varphi(i)$  by  $w_\alpha^i$ ,  $i \in \Lambda$ . Finally set  $\Omega = \Lambda^{\mathbb{N}}$ . We regard  $\Omega$  as part of the “boundary” of  $W^*(\Lambda)$ ; namely for sequences  $w^{(n)} \in W(\Lambda)$  of increasing length, and  $\alpha^{(n)} \in \mathcal{F}$ , we say  $w_{\alpha_n}^{(n)} \rightarrow \omega \in \Omega$  if  $\alpha_n \rightarrow \infty$  and  $w^{(n)}(p) = \omega(p)$  for each  $p \in \mathbb{N}$  and  $n > n(p)$ .

We next define subspaces of  $W(\Lambda)$  and  $W^*(\Lambda)$ . Let  $\Sigma = \{\varphi_1(t), \varphi_2(t), \dots\}$  be a sequence in  $W^*(\Lambda)$ . We define  $W_\Sigma^*(\Lambda)$  as consisting of words of the form  $\varphi_1(u_1)\varphi_2(u_2) \cdots \varphi_h(u_h)$ ,  $h$  arbitrary, with  $u_i \in \Lambda$  or  $u_i \in \Lambda \cup \{t\}$  respectively. It is clear how to define  $\Omega_\Sigma$ .

By analogy with  $\mathcal{F}$ -sequences we have

**Definition 2.** A function  $x(w_\alpha)$  is called a  $W^*(\Lambda)$ -sequence. Its restriction to a subspace is a  $W^*(\Lambda)$ -subsequence.

**Definition 3.** If  $x(w_\alpha)$  takes values in a metric space  $M$ . Then we say  $x(w_\alpha)$  is coherent if there exists a function  $x^* : \Omega \rightarrow M$  such that  $x(w_\alpha) \rightarrow x^*(\omega)$  uniformly as  $w_\alpha \rightarrow \omega \in \Omega$ .

The following theorem is a far reaching extension of the Hales-Jewett theorem. It includes Hindman’s theorem. See [Ca1, FK3].

**Theorem 4.** A  $W^*(\Lambda)$ -sequence with values in compact metric space has a coherent  $W^*(\Lambda)$ -subsequence.

Given a  $W(\Lambda)$ -system we will restrict to appropriate subspaces in order to achieve coherence of certain expressions. Note that the expression in (5.2) is a  $W^*(\Lambda)$ -sequence.

The next lemma shows how IP-systems occur within  $W(\Lambda)$ -systems.

**Lemma 5.** *Let  $w \in W_N(\Lambda)$ ,  $\alpha, \beta \subset \{1, 2, \dots, N\}$ ,  $\alpha < \beta$  and suppose that for  $n \in \alpha \cup \beta$ ,  $w(n) = i$ . Then for any  $W(\Lambda)$ -system  $\{T(w)\}$*

$$T(w_\alpha^j)T(w_\alpha^i)^{-1}T(w_\beta^j)T(w_\beta^i)^{-1} = T(w_{\alpha \cup \beta}^j)T(w_{\alpha \cup \beta}^i)^{-1}. \quad (14)$$

Using Theorem 4 we pass to a subsystem for which all

$$Q_{ij}(\omega) = \lim_{\substack{\alpha \rightarrow \infty \\ w \rightarrow \omega}} T(w_\alpha^j)T(w_\alpha^i)^{-1}$$

exist uniformly. On account of Lemma 5 and Proposition 4.7 these define conditional expectation operators  $E(\cdot | \mathcal{D}_{ij}(\omega))$  provided each  $i \in \Lambda$  occurs infinitely often in  $\omega$ . These  $Q_{ij}(\omega)$  represent the “ergodic averaging” operator in the  $W(\Lambda)$  context. Note that we obtain a “field” of operators. In a similar way we can form rigid extensions, and quasi-distal factor systems, except that the entire structure will vary (continuously) with  $\omega \in \Omega$ . The  $\sigma$ -algebras would not be invariant but for  $\alpha$  far out the operators

$$T_j^i(w_\alpha) = T(w_\alpha^j)T(w_\alpha^i)^{-1}$$

will move  $\mathcal{D}_{ij}(\omega)$  to  $\mathcal{D}_{ij}(\omega')$  with  $\omega'$  close to  $\omega$ . With this apparatus we can proceed as in Sect. 4.

## References

- [BB1] Brown, T.C., Buhler, J.P.: A density version of a geometric Ramsey theorem. *J. Comb. Theory Ser. A* **32** (1982) 20–34
- [Be1] Bergelson, V.: Weakly mixing PET. *Ergodic Theory Dyn. Systems* **7** (1987) 337–349
- [Ca1] Carlson, T.J.: Some unifying principles in Ramsey theory. *Discrete Math.* **68** (1988) 117–169
- [CL1] Conze, J.P., Lesigne, E.: Sur un théorème ergodique pour les mesures diagonales. Preprint
- [ET1] Erdős, P., Turán, P.: On some sequences of integers. *J. London Math. Soc.* **11** (1936) 261–264
- [Fu1] Furstenberg, H.: Ergodic behavior of diagonal measures and a theorem of Szemerédi on arithmetic progressions. *J. d'Analyse Math.* **31** (1977) 204–256
- [Fu2] Furstenberg, H.: Recurrence in Ergodic Theory and Combinatorial Number Theory. Princeton University Press (1981)
- [Fu3] Furstenberg, H.: Nonconventional ergodic theorems. In: *Symposium on Legacy of von Neumann. Proceedings of Symposia in Pure Mathematics*, Amer. Math. Soc. **50** (1990) 43–56
- [FK1] Furstenberg, H., Katznelson, Y.: An ergodic Szemerédi theorem for commuting transformations. *J. d'Analyse Math.* **34** (1978) 275–291

- [FK2] Furstenberg, H., Katznelson, Y.: An ergodic Szemerédi theorem for IP-systems and combinatorial theory. *J. d'Analyse Math.* **45** (1985) 117–168
- [FK3] Furstenberg, H., Katznelson, Y.: Idempotents in compact semigroups and Ramsey theory. *Israel J. Math.* **68** (1989) 257–270
- [FK4] Furstenberg, H., Katznelson, Y.: A density version of the Hales-Jewett theorem. *J. D'Analyse Math.* (to appear)
- [GR1] Graham, R.L., Rothschild, B.L.: Ramsey's theorem for  $n$ -parameter sets. *Trans. Amer. Math. Soc.* **159** (1971) 257–292
- [GRS1] Graham, R.L., Rothschild, B.L., Spencer, J.H.: *Ramsey theory*. Wiley, New York 1980
- [HJ1] Hales, A.W., Jewett, R.I.: Regularity and positional games. *Trans. Amer. Math. Soc.* **106** (1963) 222–229
- [Hi1] Hindman, N.: Ultrafilters and combinatorial number theory. In: Number theory Carbondale 1979. Lecture Notes in Mathematics, vol. 751. Springer, Berlin Heidelberg New York 1979, pp. 119–184
- [Sz1] Szemerédi, E.: On sets of integers containing no  $k$  elements in arithmetic progression. *Acta Arith.* **27** (1975) 199–245



# Random Schrödinger Operators

*Shinichi Kotani*

Department of Mathematics, University of Tokyo, Hongo, Tokyo 113, Japan

## 1. Introduction

As for the present subject of random Schrödinger operators, two articles had already appeared in the proceedings of the last ICM at Berkeley presenting various viewpoints (Pastur [13], Spencer [15]). Since then several important works have been done for Schrödinger operators with certain quasi periodic potentials in the one dimensional case. There does not seem to be much to add in this note to what already exists in the literature; however, we would like to mention several open problems in the case where random potentials are nearly almost periodic. We focus our attention only on the one dimensional case.

## 2. Setup of Random Schrödinger Operators

Let

$$\Omega = L^2_{\text{real}} \left( \mathbb{R}, \frac{dx}{1 + |x|^3} \right).$$

For  $q \in \Omega$  set

$$L(q) = -\frac{d^2}{dx^2} + q.$$

It is not difficult to see that  $L(q)$  defines a selfadjoint operator in  $L^2(\mathbb{R}, dx)$ . We consider a family of operators  $\{L(q); q \in \Omega\}$  under a probability measure  $P$  on  $\Omega$ . We impose the following conditions on  $P$ :

- (1) **Shift invariance:** For any  $A \in \mathfrak{B}(\Omega)$  and  $x \in \mathbb{R}$  it holds that

$$P(T_x A) = P(A)$$

where  $T_x$  is the shift operator on  $\Omega$ .

- (2) **Ergodicity:**  $T_x A = A$  a.e. with respect to  $P$  for every  $x \in \mathbb{R}$  implies  $P(A) = 0$  or 1.

**(3) Integrability:**

$$\int_{\Omega} \left\{ \int_0^1 q(x)^2 dx \right\} P(dq) < \infty.$$

There are many examples satisfying these conditions such as periodic potentials, almost periodic potentials and any stationary ergodic process with finite second moment.

Now we introduce quantities related to the operator  $L(q)$ , which are necessary for subsequent discussions. For any  $\lambda \in \mathbb{C} \setminus \mathbb{R}$ , there exist unique solutions  $f_{\pm}(x, \lambda, q)$  of

$$L(q)u = \lambda u, \text{ with } u(0) = 1 \text{ and } u \in L^2(\mathbb{R} \pm, dx),$$

respectively. Define

$$h_{\pm}(\lambda, q) = \pm \frac{d}{dx} f_{\pm}(x, \lambda, q)|_{x=0}.$$

It is of some interest to see

$$f_{\pm}(x, \lambda, q) = \exp \left( \pm \int_0^x h_{\pm}(\lambda, T_y q) dy \right). \quad (2.1)$$

The Green function of  $L(q)$  is given by

$$g_{\lambda}(x, y, q) = -(h_+(\lambda, q) + h_-(\lambda, q))^{-1} f_+(x, \lambda, q) f_-(y, \lambda, q)$$

if  $x \geq y$ . The importance of  $h_{\pm}$  in the study of almost periodic Schrödinger operators was first recognized by R. Johnson and J. Moser [3].

### 3. Fundamental Theorems and Open Problems

Define a holomorphic function  $w(\lambda)$  on  $\mathbb{C} \setminus \mathbb{R}$  by

$$w(\lambda) = -\frac{1}{2} \int_{\Omega} g_{\lambda}(0, 0, q)^{-1} P(dq)$$

for  $P \in \mathcal{P}(\Omega)$ . This function was first introduced by R. Johnson and J. Moser [3] in the almost periodic case. Without difficulty we find

$$w(\lambda) = \int_{\Omega} h_+(\lambda, q) P(dq) = \int_{\Omega} h_-(\lambda, q) P(dq).$$

The function  $w$  has a positive imaginary part on the upper half plane  $\mathbb{C}_+$  and such a holomorphic function is called a Herglotz function, which has appeared often in the spectral theory of differential operators and the classical theory of moment problems. Our  $w$  has more special properties by which we see the existence of the finite limit  $w(\xi + i0)$  for any  $\xi \in \mathbb{R}$ . We set for  $\xi \in \mathbb{R}$

$$w(\xi + i0) = -\gamma(\xi) + i\pi n(\xi).$$

Then it is known that  $\gamma(\xi)$  is non-negative and  $n(\xi)$  is continuous, non-negative and non-decreasing.  $\gamma(\xi)$  is called the Lyapunov exponent and  $n(\xi)$  is called the integrated density of states.  $\pi n(\xi)$  is called the rotation number. The readers will find a suitable explanation of  $\gamma, n$  in R. Johnson and J. Moser [3] or S. Kotani [5]. Now we can state

**Theorem** (S. Kotani [4]). *Set*

$$N = \{\xi \in \mathbb{R}; \gamma(\xi) = 0\}.$$

*Then for almost all  $q \in \Omega$  with respect to  $P$*

$$h_+(\xi + i0) = -\overline{h_-(\xi + i0)}$$

*holds for almost all  $\xi \in N$ .*

This theorem has several interesting implications, among which we see that, by the reflection principle,  $h_{\pm}(\lambda, q)$  have analytic continuation from  $\mathbb{C}_+$  to the lower half plane  $\mathbb{C}_-$  through the interior  $\mathring{N}$ . Therefore this combined with the formula (2.1) implies that  $L(q)$  has generalized Bloch solutions if  $\lambda \in \mathring{N}$ , and in particular  $L(q)$  has pure absolutely continuous spectrum on  $\mathring{N}$ . A more careful study shows

**Theorem** (K. Ishii [2], L.A. Pastur [12], S. Kotani [4]). *With probability one the absolutely continuous spectrum of  $L(q)$  coincides with  $\bar{N}^\mu := \{\xi \in \mathbb{R}; \mu(U \cap N) > 0$  for any neighbourhood  $U$  of  $\xi\}$ , where  $\mu$  denotes Lebesgue measure.*

Therefore we can tell everything about the absolutely continuous spectrum by looking at the function  $w$  alone. Then the following question naturally arises:

**P1.** *Can we determine the pure absolute continuity of the spectrum from  $w$  alone?*

The above argument shows that if  $(dn)(\mathbb{R} \setminus \mathring{N}) = 0$ , then  $L(q)$  has only absolutely continuous spectrum with probability one. However this condition is far from the necessity. In this respect, S. Kotani, M. Krishna [6] gives more recent information. It gives a more delicate sufficient condition which was essentially introduced by B.M. Levitan [7]

Once this problem was solved, the next question would be

**P2.** *Does the absolute continuity of the spectrum imply the almost periodicity of the random potentials?*

It is well known that if the spectrum of  $L(q)$  coincides with the set  $N$  and  $N$  consists of only finitely many closed intervals, then  $q$  must be a quasi periodic potential (S.P. Novikov [11], H.P. McKean and P. van Moerbecke [9]). We have also several papers treating potentials having infinitely many gaps of certain types in the spectrum (H.P. McKean and E. Trubowitz [10], V.A. Marchenko and I.V. Ostrovskii [8], L.A. Pastur and B.A. Tkachenko [14]). Among them B.M. Levitan [7] is very closely related to our purpose and is restated in our framework

in S. Kotani and M. Krishna [6]. W. Craig [1] provides a non-probabilistic treatment of this problem and gives a simpler proof of the trace formula. Although much effort has been made to clarify the relationship between the spectrum and the potentials, it seems that we have not grasped a key point yet.

## References

1. Craig, W.: The trace formula for Schrödinger operators on the line. Preprint
2. Ishii, K.: Localization of eigenstates and transport phenomena in one dimensional disordered systems. *Supp. Prog. Theor. Phys.* **53** (1973) 77–138
3. Johnson, R., Moser, J.: The rotation number of almost periodic potentials. *Comm. Math. Phys.* **84** (1982) 403–438
4. Kotani, S.: Ljapounov indices determine absolutely continuous spectrum of stationary random one-dimensional Schrödinger operators. *Proc. of Taniguchi Symp. SA. Katata 1982*, pp. 225–247
5. Kotani, S.: One dimensional Schrödinger operators and Herglotz functions. *Proc. Taniguchi Symp. PMMP. Katata 1985*, pp. 219–250
6. Kotani, S., Krishna, M.: Almost periodicity of some random potentials. *J. Funct. Anal.* **78** (1988) 390–405
7. Levitan, B.M.: On the closure of the set of finite-zone potentials. *Math. USSR Sb.* **51** (1985) 67–89
8. Marchenko, V.A., Ostrovskii, I.V.: A characterization of the spectrum of Hill's operator. *Math. USSR Sb.* **26** (1975) 493–554
9. McKean, H.P., van Moerbecke, P.: The spectrum of Hill's equation. *Invent. math.* **30** (1975) 217–274
10. McKean, H.P., Trubowitz, E.: Hill's operator and hyperelliptic function theory in the presence of infinitely many branch points. *Comm. Pure Appl. Math.* **29** (1976) 143–226
11. Novikov, S.P.: The periodic problem for the KdV equation I. *Funct. Anal. Appl.* **8** (1974) 236–246
12. Pastur, L.A.: Spectral properties of disordered systems in one-body approximation. *Comm. Math. Phys.* **88** (1980) 167–196
13. Pastur, L.A.: Spectral properties of metrically transitive operators and related problems. *Proc. of ICM Berkeley 1986*, pp. 1296–1311
14. Pastur, L.A., Tkachenko, B.A.: On the spectral theory of the one-dimensional Schrödinger operator with a limit periodic potential. *Dokl. Acad. Nauk USSR* **279** (1984) 1050–1054
15. Spencer, T.: Random and quasiperiodic Schrödinger operators. *Proc. of ICM Berkeley 1986*, pp. 1312–1318

# De Rham Cohomology of Wiener-Riemannian Manifolds

Shigeo Kusuoka

Research Institute for Mathematical Sciences, Kyoto University Kyoto 606, Japan

## 0. Introduction

In the study of Riemannian manifolds, we often observe operators in  $L^2$ -spaces. The study of infinite dimensional manifolds had started a long time ago, but we did not think of such  $L^2$ -spaces because of a lack of good Riemannian volumes. However, some people (e.g. Gross, Kuo, Eells, Elworthy, Ramer [7]) had tried to study infinite dimensional manifolds which possess a good Hilbert-Riemannian metric and a good measure which is modeled after the Wiener measure. Let us call such an infinite dimensional manifold a Wiener-Riemannian manifold.

After Malliavin introduced stochastic calculus of variation, which is called Malliavin calculus now, there was a great progress on stochastic analysis. (If we made a list of contributors, it might be very long). And this progress enables us to challenge the study of Wiener-Riemannian manifolds again. Such challenge started already by some people (e.g. Malliavin, Airault, Biesen [1, 2], Watanabe, Kazumi, Aida).

In this paper, we focus on the de Rham cohomology of Wiener-Riemannian manifolds and explain what the purpose of this topic is and what we can prove at the moment.

## 1. Preliminary Facts and Motivation

We say that a triple  $(\mu, H, B)$  is an abstract Wiener space, if

(1)  $B$  is a separable real Banach (or Fréchet) space,

(2)  $H$  is a separable real Hilbert space continuously and densely embedded in the Banach space  $B$ ,

and

(3)  $\mu$  is a probability measure in  $B$  such that

$$\int_B \exp(\sqrt{-1}_B \langle z, u \rangle_{B^*}) \mu(dz) = \exp\left(-\frac{1}{2} \|u\|_{H^{*2}}\right)$$

for any  $u \in B^* \subset H^*$ .

Here  $B^*$  and  $H^*$  are the dual space of  $B$  and  $H$  respectively. Often the dual space  $H^*$  is identified with the Hilbert space  $H$  itself. However, we will not take this convention here.

An important example of abstract Wiener spaces is an ordinary Wiener space  $(\mu_0, H_0, W^d)$ , where  $W^d = \{w \in C([0, 1]; \mathbb{R}^d); w(0) = 0\}$ ,  $H_0 = \{h \in W^d; h(t) \text{ is absolutely continuous in } t \text{ and } \int_0^1 |\dot{h}(t)|^2 dt < \infty\}$ , and  $\mu_0$  is the Wiener measure on  $W^d$ .

Now let  $E$  be a separable Hilbert space. For any bounded measurable map  $f: B \rightarrow E$ , we define a measurable map  $P_t f: B \rightarrow E$ ,  $t \geq 0$ , by

$$P_t f(z) = \int_B f(e^{-t/2}z + (I_H - e^{-t})^{1/2}w) \mu(dw), \quad z \in B.$$

Then it is easy to see that  $\{P_t\}_{t \geq 0}$  can be regarded as a strongly continuous contraction semigroup in  $L^p(B; E, d\mu)$ ,  $p \in (1, \infty)$ . Let  $\mathcal{L} = \mathcal{L}_{p, E}$  denote the infinitesimal generator of the semigroup  $\{P_t\}_{t \geq 0}$  in  $L^p(B; E, d\mu)$ . For any separable real Hilbert space  $E$ ,  $s \geq 0$  and  $p \in (1, \infty)$ , we define a Banach space  $\mathbf{D}_p^s(E)$  by

$$\mathbf{D}_p^s(E) = \text{Image}(I - \mathcal{L}_{p, E})^{-s/2} \text{ in } L^p(B; E, d\mu),$$

and

$$\|u\|_{\mathbf{D}_p^s} = \|(I - \mathcal{L}_{p, E})^{s/2} u\|_{L^p(B; E, d\mu)}, \quad u \in \mathbf{D}_p^s(E).$$

Also, we define  $\mathbf{D}_p^s(E)$ ,  $s < 0$ ,  $p \in (1, \infty)$ , to be the dual Banach space of  $\mathbf{D}_q^{-s}(E^*)$ , where  $\frac{1}{p} + \frac{1}{q} = 1$ . As usual, we identify the Banach space  $L^p(B; E, d\mu)$  with the dual space of  $L^q(B; E^*, d\mu)$ . Then we see that

$$\mathbf{D}_p^s(E) \supset \mathbf{D}_{p'}^{s'}(E) \text{ if } 1 < p \leq p' < \infty \quad \text{and} \quad -\infty < s \leq s' < \infty,$$

and

$$\|u\|_{\mathbf{D}_p^s(E)} = \|(I - \mathcal{L})^{s/2} u\|_{L^p(B; E, d\mu)}, \quad u \in L^p(B; E, d\mu),$$

if  $s \in (-\infty, 0]$  and  $p \in (1, \infty)$ .

We define  $\mathbf{D}_p^\infty(E)$ ,  $p \in (1, \infty)$  and  $\mathbf{D}_{p-}^\infty(E)$ ,  $p \in (1, \infty]$ , by

$$\mathbf{D}_p^\infty(E) = \bigcap_{s \in (0, \infty)} \mathbf{D}_p^s(E),$$

and

$$\mathbf{D}_{p-}^\infty(E) = \bigcap_{q \in (0, p)} \mathbf{D}_q^\infty(E).$$

Then they become Fréchet spaces naturally.

Also, we have the following.

**Proposition 1.** *For any separable real Hilbert space  $E$ , if  $u \in \mathbf{D}_2^1(E)$ , there is a  $Du \in \mathbf{D}_2^0(H^* \otimes E)$  such that*

$$\mu \left( \left\{ z \in B; \left\| \frac{1}{t} (u(z + th) - u(z)) - Du(z)h \right\|_E > \varepsilon \right\} \right) \rightarrow 0, \quad t \rightarrow 0,$$

for any  $h \in H$  and  $\varepsilon > 0$ .

The following is due to Meyer.

**Theorem 2.** For any separable real Hilbert space  $E$  and  $p \in (1, \infty)$ , the linear operator  $D$  from  $\mathbf{D}_2^1(E)$  into  $\mathbf{D}_2^0(H^* \otimes E)$  can be extended (or restricted) to a bounded linear operator from  $\mathbf{D}_p^s(E)$  into  $\mathbf{D}_p^{s-1}(H^* \otimes E)$  for any  $s \in \mathbb{R}$  and  $p \in (1, \infty)$ .

Therefore the dual operator  $D^*$  of  $D$  is a bounded linear operator from  $\mathbf{D}_p^s(H \otimes E)$  into  $\mathbf{D}_p^{s-1}(E)$  for any  $s \in \mathbb{R}$  and  $p \in (1, \infty)$ . Let  $\iota_H : H^* \rightarrow H$  be an isometry given by  $(\iota_H u, h)_H = {}_{H^*}\langle u, h \rangle_H$ ,  $h \in H$ . Then we have

$$\mathcal{L} = -\frac{1}{2} D^* \iota_H D.$$

The following is due to Malliavin and Watanabe. By virtue of it we can consider hypersurfaces in abstract Wiener spaces.

**Theorem 3.** Let  $F = (F_1, \dots, F_n) \in \mathbf{D}_{\infty-}^{\infty}(\mathbb{R}^n)$  and assume that

$$\det(\{(DF_i, DF_j)_{H^*}\}_{i,j=1,\dots,n})^{-1} \in \bigcap_{p \in (1, \infty)} L^p(B; d\mu).$$

Then the linear map  $S_F : \mathcal{S}(\mathbb{R}^n) \rightarrow \mathbf{D}_{\infty-}^{\infty}(\mathbb{R})$  given by  $S_F f(z) = f(F(z))$ ,  $z \in B$ , is extensible to a continuous linear map from  $\mathcal{S}'(\mathbb{R}^n)$  into  $\bigcup_{\substack{s < 0 \\ p \in (1, \infty)}} \mathbf{D}_p^s(\mathbb{R})$ .

In particular, we can think of a generalized Wiener functional  $\delta(F)$  if  $F$  satisfies the assumption of Theorem 3.

For any separable real Hilbert space  $E$ , we denote by  $E^{\otimes n}$  the completion of the algebraic tensor product  $E \otimes E \otimes \cdots \otimes E$ ,  $n \geq 1$ . Then  $E^{\otimes n}$  is a real Hilbert space and  $\{e_{k_1} \otimes e_{k_2} \otimes \cdots \otimes e_{k_n}\}_{k_1, \dots, k_n=1}^{\infty}$  is a complete orthonormal basis if  $\{e_k\}_{k=1}^{\infty}$  is a complete orthonormal basis of  $E$ . We also use the convention that  $E^{\otimes 0} = \mathbb{R}$ . For each  $n \geq 1$ , let  $A_n : E^{\otimes n} \rightarrow E^{\otimes n}$  be a bounded linear map given by  $A_n(h_1 \otimes h_2 \otimes \cdots \otimes h_n) = \frac{1}{n!} \sum_{\sigma \in S_n} \text{sgn}(\sigma) h_{\sigma(1)} \otimes h_{\sigma(2)} \otimes \cdots \otimes h_{\sigma(n)}$ ,  $h_1, \dots, h_n \in E$ . Here  $S_n$  is a set of permutations of  $\{1, \dots, n\}$ . We define a closed subspace  $\bigwedge^n E$  of  $E^{\otimes n}$  by  $\bigwedge^n E = \{x \in E^{\otimes n}; A_n(x) = x\}$ .

Now we define a continuous linear operator  $\tilde{d}_n : \mathbf{D}_{\infty-}^{\infty}(\bigwedge^n H^*) \rightarrow \mathbf{D}_{\infty-}^{\infty}(\bigwedge^{n+1} H^*)$  by  $\tilde{d}_n u = (n+1) A_{n+1} D u$ . Then we have  $d_{n+1} d_n = 0$ . We also define an inner product  $(\cdot, \cdot)_n$  on  $\mathbf{D}_{\infty-}^{\infty}(\bigwedge^n H^*)$  by

$$(u, v)_n = (n!)^{-1} \cdot \int_B (u(z), v(z))_{H^* \otimes n} \mu(dz), \quad u, v \in \mathbf{D}_{\infty-}^{\infty}(\bigwedge^n H^*).$$

Then we see that the formal adjoint operator  $\tilde{d}_n^*$  is a continuous linear operator from  $\mathbf{D}_{\infty-}^{\infty}(\bigwedge^{n+1} H^*)$  into  $\mathbf{D}_{\infty-}^{\infty}(\bigwedge^n H^*)$ .

The following is due to Shigekawa [8].

**Theorem 4.** (1)  $L_n = \tilde{d}_n^* \tilde{d}_n + \tilde{d}_{n-1} \tilde{d}_{n-1}^*$  has a natural self-adjoint extension in  $L^2(B; \bigwedge^n H^*, d\mu)$ ,  $n = 0, 1, \dots$ .

$$(2) \dim(\ker L_n) = \dim(\{u \in \mathbf{D}_{\omega-}^\infty(\bigwedge^n H^*); \tilde{d}_n u = 0\} / \{\tilde{d}_{n-1} v; v \in \mathbf{D}_{\omega-}^\infty(\bigwedge^{n-1} H^*)\}) \\ = \begin{cases} 1 & n = 0 \\ 0 & n \geq 1 \end{cases}.$$

Our main purpose is to show similar results for curved spaces.

## Hypersurfaces in Wiener Spaces

Let  $M$  be a compact Riemannian manifold of finite dimension embedded in the Euclidean space  $\mathbb{R}^d$ . Then for each  $x \in M$ , the tangent space  $T_x(M)$  is regarded as a subvector space in  $\mathbb{R}^d$ . Let  $P_x$  denote the orthogonal projection in  $\mathbb{R}^d$  onto  $T_x(M)$  for each  $x \in M$ . Let us think of the following stochastic differential equation.

$$\begin{cases} dX(t, x; w) = P_{X(t, x; w)} \circ dw(t), & t \in [0, 1] \\ X(0, x; w) = x \in M \end{cases}.$$

Here  $\{w(t); t \in [0, 1]\}$  is a  $d$ -dimensional Brownian motion. This is the Stroock's representation of Brownian motion on the manifold  $M$ .

Now let us take  $x, y \in M$  and fix them throughout. Let us think of a subset  $\tilde{M}_{xy} = \{w \in W^d; X(1, x; w) = y\}$ . This is not a closed set in  $W_d$  in general. Let  $v = v_{xy}$  be a measure on  $\tilde{M}_{xy}$  given by

$$v(dw) = \delta_y(X(1, x; w))\mu(dw).$$

Let  $N(w) = i_H DX(1, x; w)^*(DX(1, x; w)DX(1, x; w)^*)^{-1}DX(1, x; w)$ ,  $w \in W^d$ . Then  $N \in \mathbf{D}_{\omega-}^\infty(H^* \otimes H)$ . If  $w \in \tilde{M}_{xy}$ , then  $N(w)$  gives normal direction of  $\tilde{M}_{xy}$  in principle, and so  $I_H - N(w)$  is the orthogonal projection in  $H$  on to  $T_w(\tilde{M}_{xy})$  in principle.

Now we introduce an equivalence relation  $\sim$  on  $\mathbf{D}_{\omega-}^\infty(\bigwedge^n H^*)$  by the following.  $u \sim v$  if  $u(w)(h_1 - N(w)h_1, \dots, h_n - N(w)h_n) = v(w)(h_1 - N(w)h_1, \dots, h_n - N(w)h_n)$   $v$ -a.e.w for all  $h_1, \dots, h_n \in H$ . Also, let  $\tilde{\mathcal{A}}^n = \mathbf{D}_{\omega-}^\infty(\bigwedge^n H^*)/\sim$ . We can define a map  $d_n : \tilde{\mathcal{A}}_n \rightarrow \tilde{\mathcal{A}}_{n+1}$  by  $d_n(u/\sim) = (\tilde{d}_{n+1}u)/\sim$ . Then obviously  $d_{n+1}d_n = 0$ .

We can introduce an inner product on  $\tilde{\mathcal{A}}_n$  by  $(u/\sim, v/\sim)_n = (n!)^{-1} \int_{M_{xy}} ((A^n((I_H - N(w))^*)u(w), (\bigwedge^n((I_H - N(w))^*)v(w))_{H_0})^{*\otimes n} v(dw), u, v \in \tilde{\mathcal{A}}^n$ . Then we can see that the formal self adjoint map  $d_n^*$  is a map from  $\tilde{\mathcal{A}}_{n+1}$  into  $\tilde{\mathcal{A}}_n$ . So we have an operator  $L_n = d_n^* d_n + d_{n-1} d_{n-1}^*$  in the Hilbert space which is the completion of  $\tilde{\mathcal{A}}_n$ . It is not difficult to see that  $L_n$  is closable. So let us denote by the same  $L_n$  the Friedrichs extension of  $L_n$ .

Our main interest is to study  $\ker L_n$ . It is a quite natural guess that there is a relation between  $\ker L_n$  and the topological cohomology of  $\tilde{M}_{xy}$ . However, the topology of  $\tilde{M}_{xy}$  is not well-defined. So we have to modify this guess as follows.

Let us think of the following ordinary differential equation.

$$\begin{cases} \frac{d}{dt} Y(t, x; h) = P_{Y(t, x; h)} \cdot \dot{h}(t), & t \in [0, 1], \\ Y(0, x; h) = x \in M \end{cases}$$

for each  $h \in H_0$ , and let

$$\overline{M}_{xy} = \{h \in H_0; Y(1, x; h) = y\}.$$

Then  $\overline{M}_{xy}$  is a closed Hilbertian submanifold in  $H_0$ . A natural guess is that there is a natural isomorphism between  $\ker L_n$  and  $H^n(\overline{M}_{xy}; \mathbb{R})$ ,  $n$ th-topological cohomology of the space  $\overline{M}_{xy}$ .

The results which we can prove at the moment is the following.

**Theorem.** (1) *There is a natural map  $j_n : \ker L_n \rightarrow H^n(\overline{M}_{xy}; \mathbb{R})$ ,  $n \geq 0$ .*

(2)  *$j_n$ ,  $n = 0, 1, 2, \dots$ , is surjective.*

(3)  *$j_n$ ,  $n = 0$  or  $1$ , is injective.*

Let us give some remarks on this theorem. First, it is well-known that  $v_{xy}(\overline{M}_{xy}) = 0$ . So even the construction of the map  $j_n$  is not trivial. Also, let

$$\text{Path}_{xy} = \left\{ \gamma : [0, 1] \rightarrow M; \gamma(t) \text{ is absolutely continuous in } t \text{ and} \right. \\ \left. \int_0^1 |\dot{\gamma}(t)|_{T_{\gamma(t)}(M)}^2 dt < \infty \right\}.$$

Then we have a natural map  $\pi : \overline{M}_{xy} \rightarrow \text{Path}_{xy}$  given by  $\pi(h)(t) = Y(t, x; h)$ . It is easy to see that  $(\overline{M}_{xy}, \pi, \text{Path}_{xy})$  is a vector bundle. So we have a natural isomorphism between  $H^n(\overline{M}_{xy}; \mathbb{R})$  and  $H^n(\text{Path}_{xy}; \mathbb{R})$ . Since the topological cohomology of path spaces have been studied by topologists, we can use their results to study  $H^n(\overline{M}_{xy}; \mathbb{R})$ .

The strategy to show our theorem is the following. If we think of a finite dimensional compact manifold  $M$ , it is well known by the name of de Rham-Hodge-Kodaira theory that there is a natural isomorphism between  $\ker(d_n^* d_n + d_{n-1}^* d_{n-1})$  and  $H_n(M; \mathbb{R})$ , and this isomorphism is usually constructed via de Rham cohomology. So we have to define de Rham cohomology for Wiener-Riemannian manifolds. To compare  $\ker(d_n^* d_n + d_{n-1}^* d_{n-1})$  and de Rham cohomology, the hypoellipticity of the operator  $d_n^* d_n + d_{n-1}^* d_{n-1}$  plays a key role. Also, to compare de Rham cohomology and  $H^n(M; \mathbb{R})$ , we use Cech-de Rham's argument which requires the existence of partition of unity and the cohomology vanishing theorem of Poincare type.

In the rest of this paper, we explain in what sense we have the definition of de Rham cohomology, the existence of partition of unity, the cohomology vanishing theorem and the hypoellipticity.

### 3. Tools

Let  $X$  be a separable metric space and let  $\mathcal{P}(X)$  denote the set of all subsets of  $X$ .

**Definition 1.** We call a function  $\alpha : \mathcal{P}(X) \rightarrow [0, \infty)$  a *finite capacity on  $X$*  if the following are satisfied.

- (1)  $\alpha(\emptyset) = 0$ .
- (2)  $\alpha(A) \leq \alpha(B)$  for any  $A, B \in \mathcal{P}(X)$  with  $A \subset B$ .
- (3)  $\alpha(A_1 \cup A_2) \leq \alpha(A_1) + \alpha(A_2)$  for any  $A_1, A_2 \in \mathcal{P}(X)$ .
- (4) For any  $A \in \mathcal{P}(X)$ ,  $\alpha(A) = \inf\{\alpha(G); A \subset G, G \text{ is an open set in } X\}$ .
- (5) There is a sequence  $\{K_n\}_{n=1}^{\infty}$  of compact sets in  $X$  such that  $\alpha(X \setminus K_n) \rightarrow 0$  as  $n \rightarrow \infty$ .

We denote by  $\mathcal{CAP}(X)$  the set of finite capacities on  $X$ .

**Definition 2.** Let  $\alpha \in \mathcal{CAP}(X)$ .

(1) We say that an element  $A \in \mathcal{P}(X)$  is  $\alpha$ -quasi-closed, if there is a sequence  $\{K_n\}_{n=1}^{\infty}$  of compact sets in  $X$  such that  $K_n \subset A$ ,  $n \geq 1$ , and  $\alpha(A \setminus K_n) \rightarrow 0$ ,  $n \rightarrow \infty$ .

(2) We say that an element  $A \in \mathcal{P}(X)$  is  $\alpha$ -quasi-open if  $X \setminus A$  is  $\alpha$ -quasi-closed.

**Definition 3.** We say that an  $\alpha \in \text{Cap}(X)$  is *countably dominated* if  $\alpha$  satisfies the following.

$$\alpha\left(\bigcup_{n=1}^{\infty} A_n\right) \leq \sum_{n=1}^{\infty} \alpha(A_n) \text{ for any } \{A_n\}_{n=1}^{\infty} \subset \mathcal{P}(X).$$

*Example 1.* For each  $s \in (0, \infty)$  and  $p \in (1, \infty)$ , let

$$C_{s,p}(G) = \inf\{\|u\|_{\mathbf{D}_p^s(\mathbb{R})^p}; u \in \mathbf{D}_s^p(\mathbb{R}), u(z) \geq 1 \text{ } \mu\text{-a.e. } z \in G\}$$

for any open set  $G$  in  $B$ , and

$$C_{s,p}(C) = \inf\{C_{s,p}(G); G \text{ is an open set in } B \text{ and } C \subset G\}$$

for any subset  $C$  in  $B$ . Also, let us define  $C_{\infty} : \mathcal{P}(B) \rightarrow [0, 1]$  by

$$C_{\infty}(C) \stackrel{\text{def}}{=} \sum_{n=1}^{\infty} 2^{-n} \cdot C_{n,n}(C) \text{ for any subset } C \text{ in } B.$$

Then  $C_{s,p}$ ,  $C_{\infty}$  are countably dominated finite capacities in  $B$ . These capacities were first introduced by Malliavin.

*Example 2.* For any  $s \in (0, \infty)$ ,  $p \in (1, \infty)$ , we define

$$\text{Cap}_{s,p}^{\infty}(G) = \inf\{\|u\|_{\mathbf{D}_p^s(\mathbb{R})^p}; u \in \mathbf{D}_{\infty}^{\infty}(\mathbb{R}), u(z) \geq 1, \mu\text{-a.e. } z \in G\}$$

for each open set  $G$  in  $B$ , and

$$\text{Cap}_{p,s}^{\infty}(C) = \inf\{\text{Cap}_{p,s}^{\infty}(G); C \subset G, G \text{ is open}\}$$

for any subset  $C$  in  $B$ . Also, let us define  $\text{Cap}^{\infty} : \mathcal{P}(B) \rightarrow [0, 1]$  by

$$\text{Cap}^{\infty}(C) \stackrel{\text{def}}{=} \sum_{n=1}^{\infty} \text{Cap}_{n,n}^{\infty}(C) \text{ for any subset } C \text{ in } B.$$

Then  $\text{Cap}_{p,s}^{\infty}$  and  $\text{Cap}^{\infty}$  are finite capacities on  $B$ .

**Definition 4.** Let  $U$  be a  $C_{\infty}$ -quasi-open set. For any  $p \in (1, \infty]$  and a separable Hilbert space  $E$ , we say that  $f : U \rightarrow E$  belongs to  $\mathcal{D}_{p,\text{loc}}^w(U; E)$  if for any  $n \geq 2 +$

$[1/(p+2)]$  and a compact set  $K$  with  $K \subset U$ , there is  $\varphi \in \mathbf{D}_n^n(\mathbb{R})$  such that

- (1)  $0 \leq \varphi \leq 1$ ,  $\varphi(z) = 0$ ,  $z \in B \setminus U$ ,
- (2)  $\varphi f \in \mathbf{D}_{(n \wedge p)-1/n}^n(E)$ ,

and

- (3)  $C_\infty(K \setminus \{\varphi = 1\}) < 1/n$ .

Then we have the following ([4, Lemma (5.24)]).

**Theorem 5.** *Let  $U_n$ ,  $n = 1, 2, \dots$ , and  $U$  are  $\text{Cap}^\infty$ -quasi-open sets, and assume that*

$$(i) \quad \bigcup_{n=1}^{\infty} U_n \subset U,$$

and

$$(ii) \quad \text{for any compact set } K \text{ with } K \subset U, \text{Cap}^\infty\left(K \setminus \left(\bigcup_{k=1}^n U_k\right)\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Then there are  $\varphi_n \in \mathcal{D}_{\infty, \text{loc}}^w(U; \mathbb{R})$  such that

- (1)  $\varphi_n(z) = 0$ ,  $U \setminus U_n$ ,  $n \geq 1$ ,
- (2)  $0 \leq \varphi_n \leq 1$ ,
- (3)  $\varphi_n f \in \mathcal{D}_{p, \text{loc}}^w(U; E)$  for any  $f \in \mathcal{D}_{p, \text{loc}}^w(U_n; E)$ ,  $n \geq 1$ ,

and

(4) for any  $\varepsilon > 0$  and compact set  $K$  with  $K \subset U$ , there are  $m \geq 1$  and a compact set  $K' \subset U$  such that

- (i)  $\varphi_n(z) = 0$ ,  $z \in K'$ ,  $n \geq m + 1$ ,

$$(ii) \quad \sum_{n=1}^m \varphi_n(z) = 1, z \in K',$$

and

- (iii)  $\text{Cap}^\infty(K \setminus K') < \varepsilon$ .

Now let  $U$  be a  $C_\infty$ -quasi-open set in  $W^d$ . We give an equivalence relation  $\sim$  in  $\mathcal{D}_{p, \text{loc}}^w(U; \bigwedge^n H_0^*)$ ,  $n \geq 0$ ,  $p \in (1, \infty]$ , by the following:

$\alpha \sim \beta$  if

$$\alpha(w)(h_1 - N(w)h_1, \dots, h_n - N(w)h_n) = \beta(w)(h_1 - N(w), \dots, h_n - N(w)h_n)$$

v-a.e.  $w \in U$  for any  $h_1, \dots, h_n \in H_0$ . Now we define a vector space  $\mathcal{A}_p^n(U)$ ,  $n \geq 0$ ,  $p \in (1, \infty]$ , by  $\mathcal{A}_p^n(U) \stackrel{\text{def}}{=} \mathcal{D}_{p, \text{loc}}^w(U; \bigwedge^n H_0^*)/\sim$ . Then  $d_n: \mathcal{A}_p^n(U) \rightarrow \mathcal{A}_p^{n+1}(U)$  is well-defined and  $d_{n+1} d_n = 0$ .

For any  $\gamma \in C^\infty([0, 1]; M)$  and  $\varepsilon > 0$ , let  $U(\gamma, \varepsilon)$  be a  $\text{Cap}^\infty$ -quasi-open set given by

$$U(\gamma, \varepsilon) = \left\{ w \in W^d; \max_{t \in [0, 1]} \text{dis}_M(X(t, x; w), \gamma(t)) < \varepsilon \right\}.$$

Then we can prove the following.

**Theorem 6.** *There is an  $\varepsilon_0 > 0$  satisfying the following. If  $\gamma_k \in C^\infty([0, 1]; M)$ ,  $\varepsilon_k \in (0, \varepsilon_0)$ ,  $k = 1, \dots, n$ , and  $\nu_{xy}\left(\bigcap_{k=1}^n U(\gamma_k, \varepsilon_k)\right) > 0$ , then for any  $\alpha \in \mathcal{A}_p^m\left(\bigcap_{k=1}^n U(\gamma_k, \varepsilon_k)\right)$ ,*

$m \geq 1$ , with  $d_m \alpha = 0$ , there is a  $\beta \in \mathcal{A}_p^{m-1} \left( \bigcap_{k=1}^n U(\gamma_k, \varepsilon_k) \right)$  such that  $d_{m-1} \beta = \alpha$ .

For any  $\gamma \in C^\infty([0, 1]; M)$  and  $\varepsilon > 0$ , let  $\bar{U}(\gamma, \varepsilon)$  be an open set in  $H_0$  given by

$$\bar{U}(\gamma, \varepsilon) = \left\{ h \in H_0; \max_{t \in [0, 1]} \text{dis}_M(Y(t, x; h), \gamma(t)) < \varepsilon \right\}.$$

Then the support theorem tells us the following

**Proposition 7.** Let  $\gamma_k \in C^\infty([0, 1]; M)$  and  $\varepsilon_k \in (0, \varepsilon_0)$ ,  $k = 1, \dots, n$ . Then  $v_{xy} \left( \bigcap_{k=1}^n U(\gamma_k, \varepsilon_k) \right) > 0$  if and only if  $\overline{M}_{xy} \cap \left( \bigcap_{k=1}^n \bar{U}(\gamma_k, \varepsilon_k) \right) \neq \emptyset$ .

Then by using Cech-de Rham's argument, we have the following.

**Corollary 8.**  $\{\alpha \in \mathcal{A}_p^n(W^d); d_n \alpha = 0\}/\{d_{n-1} \beta; b \in \mathcal{A}_p^{n-1}(W^d)\} \simeq H^n(\overline{M}_{xy}; \mathbb{R})$  for any  $n \geq 0$  and  $p \in (1, \infty]$ .

Also, we have the following hypoelliptic result

**Theorem 9.**  $\text{Ker } L_n \subset \mathcal{A}_2^n(W^d)$ ,  $n \geq 0$ .

## References

1. Airault, H., Biesen J.: Géométrie riemannienne en codimension finie sur l'espace de Wiener. *C.R. Acad. Sci. Paris Série I*, **311** (1990) 125–130
2. Airault, H., Malliavin, P.: Intégration géométrique sur l'espace de Wiener. *Bull Sci. Math.* **112** (1988) 13–55
3. Kusuoka, S.: On the foundation of Wiener-Riemannian manifolds. In: Elworthy, K.D., Zambrini, J-C. (eds) Stochastic analysis, path integration and dynamics. (Pitman Research Notes in Math. Series vol. 200). Longman Scientific & Technical, Essex 1989, pp. 130–164
4. Kusuoka, S.: Analysis on Wiener spaces I, Nonlinear maps. (To appear in *J. Funct. Anal.*)
5. Kusuoka, S.: Analysis on Wiener spaces II, Differential forms. Preprint RIMS Kyoto Univ
6. Kusuoka, S.: Analysis on Wiener spaces III. In preparation
7. Ramer, R.: On the de Rham complex of finite codimensional differential forms on infinite dimensional manifolds. Preprint 1974
8. Shigekawa, I.: De Rham-Hodge-Kodaira's decomposition on an abstract Wiener space. *J. Math. Kyoto Univ.* **26** (1986) 191–202

# Some Recent Results in the Asymptotic Theory of Statistical Estimation

Lucien M. Le Cam\*

Department of Statistics, University of California, Berkeley, CA94720, USA

## 1. Introduction

One of the simplest results in asymptotic theory of estimation is the Hájek-Le Cam asymptotic minimax theorem. Besides being simple, it has many applications. We review the theorem and give brief indications on some applications.

The theorem is called Hájek-Le Cam because it was proved by Hájek (1972) for the asymptotically normal (more precisely LAN) case. There was a previous theorem by Le Cam (1953). Hájek's result was substantially extended in Le Cam (1979).

Section 2 below gives a summary of definitions and notation. Section 3 reviews the asymptotic minimax theorem. Section 4 indicates how the theorem can be applied to problems recently studied by Donoho and Liu (1990), by M. Low (1989) and by Golubev and Nussbaum (1990). For further applications of the asymptotic minimax theorem, see Millar (1983).

## 2. Definitions and Notation

We shall use the definitions of Le Cam (1986) with indication of conditions under which these definitions reduce to the more usual ones.

An *experiment*  $\mathcal{E} = \{P_\theta; \theta \in \Theta\}$  will be given by a  $\sigma$ -field  $\mathcal{A}$  carried by a set  $\mathcal{X}$  and a family  $\{P_\theta; \theta \in \Theta\}$  of probability measures on  $\mathcal{A}$ . The set  $\Theta$  is usually called the parameter space. The  $L$ -space  $L(\mathcal{E})$  of an experiment  $\mathcal{E}$  is the set of all finite signed measures defined on  $\mathcal{A}$  and dominated by some convergent sum  $\sum_\theta c_\theta P_\theta$ ,  $c_\theta \geq 0$ ,  $\sum_\theta c_\theta < \infty$ . Let  $\mathcal{E}$  and  $\mathcal{F}$  be two experiments, with  $\mathcal{E} = \{P_\theta; \theta \in \Theta\}$  on a  $\sigma$ -field  $\mathcal{A}$  and  $\mathcal{F} = \{Q_\theta; \theta \in \Theta\}$  on some other  $\sigma$ -field  $\mathcal{B}$ . A *transition*  $T$  from  $L(\mathcal{E})$  to  $L(\mathcal{F})$  is a positive linear map from  $L(\mathcal{E})$  to  $L(\mathcal{F})$  such that  $\|T\mu\| = \|\mu\|$  if  $\mu \geq 0$ . Here  $\|\mu\|$  is the  $L_1$ -norm  $\|\mu\| = \sup_f \{|\int f d\mu|; |f| \leq 1, f \text{ measurable}\}$ . The deficiency  $\delta(\mathcal{E}, \mathcal{F})$  is the number  $\delta(\mathcal{E}, \mathcal{F}) = \inf_T \sup_\theta \|Q_\theta - TP_\theta\|$  where the inf is over all transitions. The distance  $\Lambda(\mathcal{E}, \mathcal{F})$  is  $\max\{\delta(\mathcal{E}, \mathcal{F}), \delta(\mathcal{F}, \mathcal{E})\}$ . Two experiments  $\mathcal{E}$  and  $\mathcal{F}$  are equivalent if  $\Lambda(\mathcal{E}, \mathcal{F}) = 0$ .

---

\* Research supported by NSF grant DMS-8701426.

The reader who would prefer working only with transitions given by Markov kernels can satisfy himself or herself that all transitions from  $L(\mathcal{E})$  to  $L(\mathcal{F})$  are given by Markov kernels if 1) the family  $\{P_\theta\}$  is dominated and 2) the  $Q_\theta$  are Borel measures on a Borel subset of a complete separable metric space.

An estimation problem consists of an experiment  $\mathcal{E} = \{P_\theta; \theta \in \Theta\}$  together with a set  $Z$  and a loss function  $W$  defined on  $\Theta \times Z$  to  $(-\infty, +\infty]$  such that  $\inf_z W_\theta(z) > -\infty$ . The set  $Z$  will also be assumed to carry a vector lattice  $\Gamma'$  of bounded numerical functions, complete for the sup norm and such that  $1 \in \Gamma'$ .

A decision procedure  $\varrho$  is then a transition  $\varrho$  from  $L(\mathcal{E})$  to the dual  $\Gamma'$  of  $\Gamma$  (for the sup norm). Such a transition has a value  $\gamma\varrho P$  for  $\gamma \in \Gamma$  and  $P \in L(\mathcal{E})$ . (This is a contraction of  $\int [\int \gamma(z)K(dz, x)]P(dx)$ .) The risk of  $\varrho$  at  $\theta$  is  $R(\theta, \varrho) = W_\theta\varrho P_\theta = \sup\{\gamma\varrho P_\theta; \gamma \in \Gamma, \gamma \leq W_\theta\}$ .

Here again the reader who prefers to work with Markov kernels ( $K$ , as above) can assume that 1) the  $\{P_\theta\}$  are dominated 2)  $Z$  is compact,  $\Gamma = C(Z)$  and each  $W_\theta$  is lower semicontinuous.

An estimation problem given by an experiment  $\mathcal{E} = \{P_\theta; \theta \in \Theta\}$  and a loss function  $W$  has a set  $\mathcal{R}(\mathcal{E}, W)$  of possible risk functions, the set of functions  $f$  from  $\Theta$  to  $(-\infty, +\infty]$  such that there is a decision procedure  $\varrho$  for which  $W_\theta\varrho P_\theta \leq f(\theta)$  for all  $\theta \in \Theta$ .

Often we shall need to work with subsets  $F \subset \Theta$ . Then  $\mathcal{E}_F$  will be  $\mathcal{E}_F = \{P_\theta; \theta \in F\}$ .

### 3. The Asymptotic Minimax Theorem

The distance defined in Sect. 2 gives a topology on the set of (equivalence classes) of experiments indexed by a set  $\Theta$ . Another topology is the weak topology: A directed set  $\{\mathcal{E}_v\}; \mathcal{E}_v = \{P_{\theta,v}; \theta \in \Theta\}$  converges weakly to  $\mathcal{F}$  if for every finite set  $F \subset \Theta$ , the distances  $\Delta(\mathcal{E}_{v,F}, \mathcal{F}_F)$  tend to zero. This is equivalent to convergence in distribution of the vector of likelihood ratios  $\left\{ \frac{dP_{t,v}}{dP_{s,v}}, t \in F \right\}$  for all  $s \in F$ .

To state the theorem call a loss function  $V$  special if  $V_\theta \in \Gamma$  for each  $\theta \in \Theta$ .

**Theorem 1.** Let  $f$  be a function that does not belong to  $\mathcal{R}(\mathcal{F}, W)$ . Then there is a special  $V \leq W$ , a number  $\alpha > 0$ , a finite set  $F$  and an  $\varepsilon > 0$  such that if  $\Delta(\mathcal{E}_F, \mathcal{F}_F) < \varepsilon$  then  $f + \alpha$  restricted to  $F$  does not belong to  $\mathcal{R}(\mathcal{E}_F, V)$ .

For a proof, see Le Cam (1979) or Le Cam (1986) pp. 109–110.

**Remark 1.** There is a weaker version of the theorem that might be easier to visualize. Let  $\{\mathcal{E}_v\}$  be a directed family of experiments  $\mathcal{E}_v = \{P_{\theta,v}; \theta \in \Theta\}$ . Assume that the  $\mathcal{E}_v$  converge weakly to  $\mathcal{F}$  and that for each  $v$  the function  $f_v$  belongs to  $\mathcal{R}(\mathcal{E}_v, W)$ . If  $f_v$  converges pointwise to  $f$  then  $f \in \mathcal{R}(\mathcal{F}, W)$ . This easy version is not sufficient for applications where one wants to truncate  $W$ . The fact that the finite set  $F$  the special  $V$  and the  $\varepsilon$  of Theorem 1 depends only on the triplet  $(\mathcal{F}, W, f)$  is also lost in the weaker version.

*Remark 2.* Theorem 1 has been stated in the general framework of Sect. 2 with procedures that are “transitions”. If one wants to restrict oneself to transitions representable by Markov kernels it is sufficient to put restrictions on the limit  $\mathcal{F}$  and the loss  $W$ . Call  $\mathcal{R}(\mathcal{F}, W, \text{Markov})$  the set of functions defined as in Sect. 2 for  $\mathcal{R}(\mathcal{F}, W)$  but for transitions that are Markov kernels. It is enough to assume that  $\mathcal{R}(\mathcal{F}, W) = \mathcal{R}(\mathcal{F}, W, \text{Markov})$  for the limit experiment  $\mathcal{F}$ . Assumptions that insure this are given in Sect. 2 and in Le Cam (1986) pp. 11–14. No assumptions need to be placed on the experiments  $\mathcal{E}$  such that  $A(\mathcal{E}, \mathcal{F}) < \varepsilon$ .

Theorem 1 uses only weak convergence to  $\mathcal{F}$  of the experiments  $\mathcal{E}$ . There is another mode of convergence that is usually available at very little cost. It is as follows.

Take a fixed  $\mathcal{F} = \{Q_\theta; \theta \in \Theta\}$  and call a set  $S \subset \Theta$  compact if the set  $\{Q_\theta; \theta \in S\}$  is compact in  $L(\mathcal{F})$  for the  $L_1$ -norm. Let  $\{\mathcal{E}_v\}$  be a directed family of experiments,  $\mathcal{E}_v = \{P_{\theta,v}; \theta \in \Theta\}$ . It is said to converge to  $\mathcal{F}$  on compacts if for each compact set  $S$  the restrictions  $\mathcal{E}_{v,S}$  are such that  $A(\mathcal{E}_{v,S}, \mathcal{F}_S)$  tends to zero.

The standard LAN conditions of Le Cam (1960) imply convergence on compacts. (Hájek's 1972 do not.) According to Lindae (1972) convergence on compacts follows from pointwise convergence plus some tail equicontinuity of differences  $\|P_{s,v} - P_{t,v}\|$ ,  $s, t \in S$  compact. In many cases one would wish to consider convergence on precompact sets instead of compacts. The precompact convergence can be reduced to the compact one by completing  $\mathcal{F}$ . This can be achieved without any difficulty.

Now if  $\mathcal{E}_v$  converges on compacts to  $\mathcal{F}$ , Theorem 1 is certainly applicable, but can one say more? In the direction of lower bounds for the risk, perhaps very little can be said. However here are two results, that are of some interest.

**Theorem 2.** Assume that, for compacts defined as above,  $\mathcal{E}_v$  converges to  $\mathcal{F}$  on compacts and that  $W$  is bounded (that is  $\sup\{|W_\theta(z)|; \theta \in \Theta, z \in Z\} < \infty$ .) Then if  $f \in \mathcal{R}(\mathcal{F}, W)$  there is for each  $v$  an  $f_v \in \mathcal{R}(\mathcal{E}_v, W)$  such that  $f_v \rightarrow f$  uniformly on the compact subsets of  $\Theta$ .

This is easy to see. It tends to indicate that some results that can be achieved on the limit  $\mathcal{F}$  can also be achieved asymptotically on the directed set  $\{\mathcal{E}_v\}$ .

Another result extends the lower bound of Theorem 1. To state it, let  $W_\theta^c = c \wedge W_\theta$  for  $c \geq 0$ . For risk functions  $W_\theta \sigma_v P_{\theta,v}$  that might not be measurable, let  $\int_* W_\theta \sigma_v P_{\theta,v} \mu(d\theta)$  be the lower integral, supremum of integrals of measurable functions not exceeding  $W_\theta \sigma_v P_{\theta,v}$ . Consider an experiment  $\mathcal{F} = \{Q_\theta; \theta \in \Theta\}$  and loss functions satisfying the following assumption:

- (A) If  $\Theta$  is pseudometrized by the distance  $d(s, t) = \|Q_s - Q_t\|$  then the risk functions  $W_\theta^c Q_\theta$  are Borel measurable in  $\theta$  for all  $c$  and all procedures  $Q$  available in  $\mathcal{F}$ .

We shall state our next theorem assuming that  $d(s, t) = \|Q_s - Q_t\|$  is in fact a metric on  $\Theta$ . Modifications for a more general case are easy.

**Theorem 3.** Suppose that condition (A) above is satisfied for the experiment  $\mathcal{F}$  and that  $d(s, t)$  defined above is a metric.

Let  $\mu$  be a finite Radon measure on  $\Theta$  (metrized by  $d$ ). Assume that  $W \geq 0$ , and let  $A = \inf_{\varrho} \int W_{\theta} \varrho Q_{\theta} \mu(d\theta)$  be the Bayes risk for  $\mu$  and  $\mathcal{F}$ .

Then for each  $b < A$  there is a  $c < \infty$ , a compact  $K \subset \Theta$  and an  $\alpha > 0$  such that if  $\Delta(\mathcal{E}_K, \mathcal{F}_K) < \alpha$  then  $\inf_{\sigma} \int_* I_K(\theta) W_{\theta}^c \sigma P_{\theta} \mu(d\theta) \geq b$ , the infimum being over all procedures  $\sigma$  available for  $\mathcal{E} = \{P_{\theta}; \theta \in \Theta\}$ .

*Proof.* Let  $\varrho_0$  be a Bayes procedure for  $\mathcal{F}$ ,  $W$  and  $\mu$ . Let  $\mathcal{V}$  the class of special loss functions  $V$  with  $V \leq W$ . Then, by definition,  $W_{\theta} \varrho_0 Q_{\theta} = \sup_{V \in \mathcal{V}} V_{\theta} \varrho_0 Q_{\theta} = \sup_c \sup_{\mathcal{V}} V_{\theta}^c \varrho_0 Q_{\theta} = \sup_c W_{\theta}^c \varrho_0 Q_{\theta}$ . Thus if the  $W_{\theta}^c \varrho_0 Q_{\theta}$  are measurable  $\int W_{\theta} \varrho_0 Q_{\theta} \mu(d\theta) = \sup_c \int W_{\theta}^c \varrho_0 Q_{\theta} \mu(d\theta)$ . Since  $W \geq 0$ , this implies that for any number  $b'$ ,  $b < b' < A$  there is a finite  $c$  and a compact  $K \subset \Theta$  such that  $\int_K W_{\theta}^c \varrho_0 Q_{\theta} \mu(d\theta) > b'$ . Let  $\alpha > 0$  be such that  $b + \|\mu\|c\alpha < b'$ . If  $\Delta(\mathcal{E}_K, \mathcal{F}_K) < \alpha/2$  there is a transition  $T$  from  $L(\mathcal{F}_K)$  to  $L(\mathcal{E}_K)$  such that  $\|P_{\theta} - TQ_{\theta}\| < \alpha$  for all  $\theta \in K$ . This  $T$  extends to a transition from  $L(\mathcal{F})$  to  $L(\mathcal{E})$ . Thus, if  $\sigma$  is any procedure on  $\mathcal{E}$ , the procedure  $\varrho = \sigma T$  defined for  $\mathcal{F}$  is such that  $|W_{\theta}^c(\sigma T)Q_{\theta} - W_{\theta}^c \sigma P_{\theta}| < c\alpha$  for all  $\theta \in K$ . This implies

$$\begin{aligned} \int_{*} (W_{\theta} \sigma P_{\theta}) \mu d(\theta) &\geq \int_{*} I_K(\theta) W_{\theta} \sigma P_{\theta} \mu(d\theta) \\ &\geq \int_{*} I_K(\theta) W_{\theta}^c \sigma P_{\theta} \mu(d\theta) \\ &\geq \int_K W_{\theta}^c(\sigma T)Q_{\theta} \mu(d\theta) - \|\mu\|c\alpha \geq b. \end{aligned}$$

Hence the result.  $\square$

**Remark 1.** It should be noted that the measurability requirement (A) is imposed only on the limit experiment  $\mathcal{F}$ , not on the approximating experiments  $\mathcal{E}$ . In the cases considered in the literature the functions  $\theta \mapsto W_{\theta}^c \varrho Q_{\theta}$  are in fact continuous. Thus measurability is not a serious problem. However it seems to be needed for the validity of Theorem 3.

**Remark 2.** Let  $\mathcal{M}$  be a class of Radon probability measures on  $\Theta$ . The conclusion of the theorem can be replaced by: Let  $b$  denote any number strictly inferior to  $\sup_{\mu} \inf_{\varrho} \int W_{\theta} \varrho Q_{\theta} \mu(d\theta)$ . Then there is a compact  $K \subset \Theta$  and numbers  $\alpha > 0$  and  $c < \infty$  such that if  $\Delta(\mathcal{E}_K, \mathcal{F}_K) < \alpha$  one has

$$\sup_{\mu} \inf_{\sigma} \int_{*} W_{\theta}^c \sigma P_{\theta} \mu(d\theta) \geq b.$$

This can be seen as in Theorem 3 taking a Bayes procedure  $\varrho_0$  for a  $\mu$  that almost achieves the  $\sup_{\mu}$  for procedures on  $\mathcal{F}$ .

**Remark 3.** One might ask whether the conclusion of Theorem 3 would remain valid under only weak convergence of the experiments instead of compact convergence. This is perhaps not so. The difficulty arises from the fact that pointwise convergence of a bounded directed set of functions does not imply convergence of their integrals.

## 4. Some Applications

A) Let us start by an example of M. Low (1989) since it is very simple. Consider, on the line  $\mathbb{R}$ , a fixed probability density  $f_0$  (with respect to Lebesgue measure) such that  $f_0(0) > 0$ ,  $\sup_x f(x) < \infty$  and such that  $f_0$  be continuous at zero. Let  $\{\alpha_n\}$  and  $\{\beta_n\}$  be nondecreasing sequences of positive numbers such that  $\alpha_n \rightarrow \infty$  and  $(\alpha_n^2 \beta_n)(f_0(0)n)^{-1} \rightarrow 1$ . Consider the class  $H$  of functions from  $\mathbb{R}$  to  $\mathbb{R}$  such that  $\int h^2 < \infty$ ,  $\int |h| < \infty$  and  $\sup_x |h(x)| < \infty$ . Let  $h_n$  be the number  $h_n = \int \alpha_n^{-1} h(\beta_n x) f_0(x) dx$ . Define  $f_n(h, x) = [1 + \alpha_n^{-1} h(\beta_n x) - h_n] f_0$  if  $1 + \alpha_n^{-1} h(\beta_n x) - h_n \geq 0$ . Let  $f_n(h, x) = f_0(x)$  otherwise. The standard Gaussian shift experiment  $\mathcal{G}$  of  $H$  is one where one takes under  $\theta = 0$  the distribution  $G_0$  of a Gaussian linear process  $Z$  indexed by  $H$  and such that  $E\langle Z, h \rangle = 0$  and  $E|\langle Z, h \rangle|^2 = \|h\|^2 = \int h^2(x) dx$ . For another value  $h \in H$  one takes for  $G_h$  the measure  $dG_h = \exp\{\langle Z, h \rangle - \frac{1}{2}\|h\|^2\} dG_0$ .

Now let  $\mathcal{E}_n = \{P_h^n; h \in H\}$  be defined by taking for  $P_h^n$  the joint distribution of  $n$  independent observations from the density  $f_n(h, x)$ . Low shows that  $\mathcal{E}_n$  converges weakly to the Gaussian  $\mathcal{G}$  as  $n \rightarrow \infty$ .

By restricting oneself to subsets of  $H$  one can obtain a variety of results from Theorem 1 (or 2). For instance Low considers a set of densities subject to a condition  $\sup_x |f^k(x)| \leq M$  and estimates of  $f(0)$ . By selecting  $\alpha_n = c_1 n^{k(2k+1)^{-1}}$  he shows that the appropriate rate of convergence of the estimate is in  $n^{k(2k+1)^{-1}}$ . This was known otherwise but Low obtains the exact limit of the risk for several loss functions.

The technique of rescaling through coefficients  $\alpha_n$  and  $\beta_n$  had been previously used by Has'minskii (1979) to study estimation of a mode. For  $\beta_n \equiv 1$ , it has been used extensively.

B) A more complicated example appears in a paper by Golubev and Nussbaum (1990). They consider the problem of estimating a signal  $t \sim f(t)$ ,  $t \in [0, 1]$  when the observations are of the form  $Y_i = f(x_{i,n}) + \xi_i$ ,  $i = 1, \dots, n$  with, for instance  $x_{i,n} = i/n$  and with noise  $\xi$  where the  $\xi_i$  are independent, mean zero, fixed variance  $\sigma^2$  and fourth moment  $E\xi_i^4$  less than a fixed constant  $c$ . The problem has been studied by many authors. A major breakthrough is due to Pinsker (1980) who considered the case where the  $\xi_i$  are Gaussian. Pinsker and subsequent authors consider the Sobolev class  $W_2^m = \{f \in L_2; D^m f \in L_2\}$  where  $L_2$  is the Hilbert space of the Lebesgue measure on  $[0, 1]$ . For the subset  $W_2^m(B) = \{f \in W_2^m; \|D^m f\|^2 \leq B\}$  let  $\Delta = \lim_n \inf_{\hat{f}} \sup_f n^{(2m)(2m+1)^{-1}} E_{f,n} \|\hat{f} - f\|^2$  where the sup is on  $f \in W_2^m(B)$  and the inf is over all estimators depending on  $n$  observations. The papers of Pinsker (1980) and Nussbaum (1985) give the result

$$\Delta = \gamma(m) B^r \sigma^{4mr} \text{ for } r = (2m+1)^{-1}$$

and  $\gamma(m) = (2m+1)^r [m/\pi(m+1)]^{2mr}$ . The fact that the  $\xi_i$  were Gaussian was essential in the proofs. Golubev and Nussbaum use only the restrictions  $E\xi_i = 0$ ,  $E\xi_i^2 = \sigma^2$ ,  $E\xi_i^4 \leq c$  and obtain a similar result.

The proof is full of ingenious devices. The relation with Theorems 1, 2 and 3 is obtainable through a series of arguments that go about as follows. Consider

a particular  $f_0$ , for instance  $f_0 \equiv 0$  and deviations from it. Let  $W_2^{m,0}$  on  $[0, 1]$  be that part of  $W_2^m$  formed by functions whose derivatives of order  $0, 1, \dots, m$  vanish at 0 and 1. For  $f \in W_2^{m,0}$  one can obtain an orthogonal expansion  $f = \sum_j c_j \varphi_j$  with  $\|\varphi_j\| = 1$  and  $\|D^m \varphi_j\|^2 = \lambda_j$  increasing in  $j$ . Now take an integer  $q$  and for  $k = 1, 2, \dots, n$  let  $I_{k,q} = ((k-1)/q, k/q]$ . Transport  $W_2^{m,0}$  to  $I_{k,q}$ , by proper scaling. Look at deviations of the type  $\sum_k \sum_{j=1}^s \varphi_{j,k,q}(x_i) f_{j,k}$  where  $\varphi_{j,k,q}$  is  $\varphi_j$  transported to  $I_{k,q}$  and put equal to zero outside  $I_{k,q}$ . Take only deviations that remain in  $W_2^m(B)$ . This allows to separate the observations by classes, the  $k$ -th class yielding a model  $y_i = \sum_{j=1}^s \varphi_{j,k,q}(x_{i,n}) f_{j,k} + \xi_i$  for those  $x_{i,n}$  that fall in  $I_{k,q}$ .

Golubev and Nussbaum let  $q$  depend on  $n$ , so it becomes  $q(n)$  of the order of  $n^r$ . They then proceed to show that the part of the regression model restricted to one of the intervals  $I_{k,q}$  converges to a Gaussian shift one.

Selecting the parameters  $f_{j,k}$  independently according to some measure  $v$  one can try to find a lower bound on the Bayes risk.

The bound in the limit is given by Theorem 1 or 3 for each of the subintervals  $I_{k,q}$ ;  $k = 1, 2, \dots, q$ . Since the Bayes risk for the entire problem is  $q(n)$  times the risk on each  $I_{k,q}$  the global lower bound can be computed for each fixed  $s$ . Then one will let  $s$  tend to infinity. Of course this is only a brief sketch of the method of proof. There are many other difficult steps on the way. One of them is to make sure that the product measure  $v^{sq}$  on  $\mathbb{R}^{sq}$  concentrates on the Sobolev ball  $W_2^m(B)$ . This was also crucial in Pinsker (1980).

In Low (1989) or Golubev and Nussbaum (1990) Theorems such as Theorems 1, 2 and 3 are used to reduce a complex problem to one in which the distributions are Gaussian and where one can often get more precise information.

C) The estimation problem treated by Donoho and Liu (1990) differs considerably from the one described in (B) above. Yet the two are closely connected. Let  $\mathcal{F}$  be a class of probability densities with respect to Lebesgue measure  $\lambda$  on an interval  $[-a, +a]$  of the line. Assume that  $\mathcal{F}$  is convex, closed and bounded for the  $L_2$ -norm,  $\|f\|^2 = \int f^2 d\lambda$ . Donoho and Liu study the problem of estimating the value  $T(f)$  of a real valued *linear* function  $T$  defined on  $\mathcal{F}$  when one takes  $n$  independent observations  $X_1, \dots, X_n$  from some  $f \in \mathcal{F}$ . For example one may want to estimate the value at zero of the  $k$ -th derivative of  $f$  subject to a local constraint on the  $m$ -th derivative, with  $k \leq m$ .

Let  $v_n$  be the empirical measure of the first  $n$  observations. One can either limit oneself to estimates  $\hat{T}$  that are *linear affine* in  $v_n$  (with risk  $R_A$  indicated by a suffix A) or use any arbitrary measurable function  $\hat{T}$  of  $v_n$  (with risk  $R_M$ , indicated by a suffix M). A first remark is that, for *affine* estimates and square loss the problem of estimation of  $T$  is not more difficult than the estimation problem for a certain Gaussian shift experiment where one observes  $Y = f + \sigma_n W$ ,  $f \in \mathcal{F}$ ,  $W$  a white noise or a Gaussian process defined on subsets of  $[-a, +a]$ , with expectations zero and a given covariance function. This is quite analogous to (B) above, but now we need to estimate only the value of  $T(f)$  instead of the whole  $f$  as in (B).

Let  $\mathcal{G}_n$  be the Gaussian experiment with observations  $Y = f + \sigma_n W$ ,  $f \in \mathcal{F}$ . Donoho and Liu proceed as follows

1)  $\mathcal{F}$  being as described, there is a worst pair  $(f_{0,n}, f_{1,n})$  of elements of  $\mathcal{F}$  such that the minimax risk for *affine* estimates and for the one dimensional system  $S_n = \{f_{\theta,n} = (1 - \theta)f_{0,n} + \theta f_{1,n}; \theta \in [0, 1]\}$  is the same as the minimax risk for affine estimates for the entire  $\mathcal{G}_n$ . Furthermore the estimate for the worst pair is minimax for  $\mathcal{G}_n$  among affine estimates. It is given by an explicit formula.

2) Consider the problem of estimating  $\theta$  for the segment  $S_n$  described above and observations

$$\int u(t)Y(dt) \text{ where } u = (f_{1,n} - f_{0,n})\|f_{1,n} - f_{0,n}\|^{-1}.$$

By sufficiency, this is equivalent to the problem where all of  $Y$  would be observed.

For the problem the risk  $R_A$  for affine estimates is a certain function  $\sigma \sim R_A(\sigma)$  of the standard deviation  $\sigma$  of  $\int u(t)Y(dt)$ . Similarly for the minimax risk  $R_M(\sigma)$  for all measurable estimates. From Ibragimov-Has'minskii (1984) one knows that  $\sup_\sigma R_A(\sigma)/R_M(\sigma)$  is bounded by a constant  $\mu^*$ . From Donoho, Liu and McGibbon (1989) one knows that  $\mu^* \leq 5/4$ . This essentially solves the problem for the Gaussian case, at least if one considers a 25% margin acceptable.

The method "almost" solves the initial problem of estimation of  $T$  defined on  $\mathcal{F}$  for the independent observations  $X_1, \dots, X_n$ , at least if one selects  $\sigma_n$  and the white noise  $W$  properly, since for *affine* estimates the two problems are essentially asymptotically equivalent. (Asymptotically only because to get exact equivalence one has to select the Gaussian set function  $W$  with a covariance that depends on the true  $f_0$ ). However that is for *affine* estimates. Would there be a possibility of doing much better for estimation of  $T(f)$  by general measurable functions of the  $X_1, \dots, X_n$ ?

Donoho and Liu resolve the difficulty, at least for usual cases, by an appeal to a theorem similar to Theorem 2, Sect. 3 above.

Let  $P_{\theta,n}$  be the joint distribution of  $X_1, \dots, X_n$  for the densities  $f_{\theta,n} = (1 - \theta)f_{0,n} + \theta f_{1,n}$ ,  $\theta \in [0, 1]$ . Consider the experiments  $\mathcal{E}_n = \{P_{\theta,n}; \theta \in [0, 1]\}$ . Consider also a Gaussian experiment

$$\mathcal{F}_n = \{Q_{\theta,n}; \theta \in [0, 1]\}$$

where  $Q_{\theta,n}$  is  $\mathcal{N}(\theta, \sigma_n^2)$  on the line. One can prove the following

**Proposition 1.** Assume that the Lévy distance between the distribution under  $P_{0,n}$  of  $\sum_{j=1}^n \left[ \frac{f_{1,n}(X_j)}{f_{0,n}(X_j)} - 1 \right]$  and a normal distribution  $\mathcal{N}(0, \tau_n^2)$  tends to zero as  $n \rightarrow \infty$ . Assume that  $\tau_n$  stays bounded. Then if  $\tau_n^2 \sigma_n^2 \rightarrow 1$  the distance  $\Delta(\mathcal{E}_n, \mathcal{F}_n)$  between the experiments  $\mathcal{E}_n = \{P_{\theta,n}; \theta \in [0, 1]\}$  and the Gaussian  $\mathcal{F}_n$  tends to zero.

This is easy to see. It follows then that the difference between the minimax risk  $R_M(\mathcal{E}_n)$  for  $\mathcal{E}_n$  and  $R_M(\mathcal{F}_n)$  for  $\mathcal{F}_n$  tends to zero.

Of course, the bulk of the argumentation of Donoho and Liu takes place on the Gaussian experiment. Donoho and Nussbaum have now extended these arguments to the estimation of certain *quadratic* functionals of the density  $f$  instead of linear ones. That the problem can be very different can be seen from an article of Bickel and Ritov (1990). The subject is still progressing.

## References

- [1] Bickel, P.J., Ritov, Y. (1990): Achieving information bounds in non and semiparametric models. *Ann. Statist.* **18**, no. 2, 925–938
- [2] Donoho, D.L., Liu, R.C. (1990): Geometrizing rates of convergence, III. Tech. Report no. 138, Department of Statistics, Berkeley
- [3] Donoho, D.L., Liu, R.C., MacGibbon, B. (1988). Minimax risk for hyperrectangles. Tech. Report no. 123, Department of Statistics, Berkeley
- [4] Donoho, D.L., Nussbaum, M. (1990): Minimax quadratic estimation of a quadratic functional. Tech. Report no. 236, Department of Statistics, Berkeley
- [5] Golubev, G.K., Nussbaum, M. (1990): A risk bound in Sobolev class regression. *Ann. Statist.* **18**, no. 2, 758–778
- [6] Hájek, J. (1972): Local asymptotic minimax and admissibility in estimation. Proc. 6th Berkeley Symp. Math. Stat. Proba., vol. 1, pp. 175–194
- [7] Has'minskii, R.Z. (1979): Lower bound for the risk of nonparametric estimates of the mode. Contribution to Statistics: J. Hájek memorial volume (Jurečková, J., ed.). Akademia, Praha, pp. 90–97
- [8] Ibragimov, I.A., Has'minskii, R.Z. (1984): On nonparametric estimation of the value of a functional in Gaussian white noise. *Theory Probab. Appl.* **24**, 18–32
- [9] Le Cam, L. (1953): On some asymptotic properties of maximum likelihood estimates and related Bayes estimates. *Univ. California Pub. Statist.* **1**, 125–142
- [10] Le Cam, L. (1960): Locally asymptotically normal families of distributions. *Univ. California Pub. Statist.* **3**, 37–98
- [11] Le Cam, L. (1979): On a theorem of J. Hájek. Contributions to statistics: J. Hájek memorial volume (Jurečková, J., ed.) Akademia, Praha, pp. 119–137
- [12] Le Cam, L. (1986): Asymptotic methods in statistical decision theory. Springer, New York Berlin Heidelberg
- [13] Lindae, D. (1972): Distributions of likelihood ratios and convergence of experiments. Unpublished Ph.D Thesis, Univ. California, Berkeley
- [14] Low, M.G. (1989): Local convergence of nonparametric density estimation problems to Gaussian shift experiments on a Hilbert space. Tech. Report no. 225, Department of Statistics, Berkeley
- [15] Millar, P.W. (1983): The minimax principle in asymptotic statistical theory. Ecole d'Eté de Probabilités, Saint Flour, 1981. (Lecture Notes in Mathematics, vol. 976.) Springer, Berlin Heidelberg New York
- [16] Nussbaum, M. (1985): Spline smoothing in regression models and asymptotic efficiency in  $L_2$ . *Ann. Statist.* **13**, 984–997
- [17] Pinsker, M.S. (1980): Optimal filtering of square integrable signals in Gaussian white noise. *Probl. Inform. Transmission* **16**, 120–133

# Localization and Intermittency: New Results

*Stanislav A. Molchanov*

Department of Mathematics and Mechanics, Moscow State University,  
Moscow, 119808, USSR

The theory of the Anderson localization was developed for a long time in the initial probabilistic framework for the operators of the form

$$H = \Delta + V(x, w), \quad x \in R^d(\mathbb{Z}^d), \quad w \in (\Omega, F, P). \quad (1)$$

Here  $\Delta$  is the Laplacian (continuous or on lattice) and  $V(x, w)$  is the random homogeneous field (or process) on a probability space  $(\Omega, F, P)$ . The central achievement of one-dimensional theory was a series of S. Kotani's articles [1]–[3], where he discovered deep connections between Ljapunov's exponents  $\gamma(\lambda)$ , the structure of the prediction of the potential  $V(\cdot)$  and the spectral localization of  $H$ .

Techniques of the cluster resolvent expansions, developed by Fröhlich and Spencer [4], together with generalization of Kotani's idea [2], proposed almost simultaneously (in different forms by Souillard et al. [5] and Simon-Wolff [6]) made it possible to solve the point spectrum problem for Anderson's tight-binding model (1) in the multidimensional lattice case  $\mathbb{Z}^d$ ,  $d > 1$ .

These results are summed up in the review article of Martinelli-Scoppola [7]. Simplifications and generalizations were proposed recently in this field by Dreifus and Klein [8].

Probabilistic approach, nevertheless, did not clear up the central physical idea of localization—the absence of resonance between a quantum particle with a given admissible energy  $\lambda$  and some (rich enough) family of the blocks of the potential. It's natural to attempt to formulate the direct geometric conditions for the individual potential  $V(\cdot)$ , which will lead to the localization of the spectrum. The first and a very important step in this direction was made in the paper Simon-Spencer [9]. They proved (in the one-dimensional case) that unboundness of potential  $V$  gives us the singular spectrum.

Moreover, in a multidimensional situation ( $d > 1$ ), even in the lattice case, there are many open problems on the thin structure of the spectrum  $H$ , and on the existence of bifurcations with respect to some parameters of the model, etc. It's related, for example, to the simplest functional of  $\sigma(H)$ , namely the integral density of states  $N(\lambda)$  (the so-called problem of Lifshitz tails).

In this lecture I'll give an account of some recent results in this direction. It is based on the papers, which I wrote together with my friends and colleagues, especially J. Görtner, A. Gordon, L. Pastur and B. Simon.

This lecture was prepared at the time of my visit to Caltech (spring 1990) and I am very grateful to Barry Simon for his hospitality, to A. Klein, R. Carmona,

B. Simon and T. Wolff for their useful discussions and to M. D'Elia, V. Jaksic, S. Katok and J. Madow for their support and assistance.

## § 1. One-Dimensional Localization. Lattice Case

The next group of results can be obtained by the combination of the ideas of Simon-Spencer [9] (appearance of the family of non-resonant blocks), Fröhlich-Spencer [4, 7] (cluster expansions of resolvent) and some new considerations, the principal of which is a randomization of energy.

**Theorem 1.** Consider in  $l^2(\mathbb{Z}_+^1)$ ,  $\mathbb{Z}_+^1 = \{0, 1, \dots\}$  the Schrödinger operator  $H^\theta = \Delta + V(x)$ ,  $x \in \mathbb{Z}_+^1$ .  $\Delta$  is a discrete Laplacian,  $\Delta\psi(x) = \psi(x+1) + \psi(x-1)$  and the parameter  $\theta$  is a boundary phase connected with the boundary condition  $\psi(-1) \cos \theta + \psi(0) \sin \theta = 0$ . Assume, that for some energy interval  $I \subset \mathbb{R}^1$  there exists the sequence of the "nonresonant blocks"  $[x_n, y_n]$ ,  $n = 1, 2, \dots$ , that is a sequence of points  $0 < x_1 \leq y_1 < x_2 \leq y_2 \dots$  such that for every  $\lambda \in I$

$$|R_\lambda^{(n)}(x_n, y_n)| \leq \delta_n ; \quad \delta_n \rightarrow 0, \quad n \rightarrow \infty. \quad (2)$$

Here  $R_\lambda^{(n)} = (H^{(n)} - \lambda)^{-1}$  is the resolvent of the operators on the block  $[x_n, y_n]$ :  $H_\psi^{(n)} = \Delta\psi + V(x)\psi(x)$ ,  $x_n \leq x \leq y_n$ ,  $\psi(x_n - 1) = \psi(y_n + 1) = 0$ .

Suppose that there exists a nondecreasing sequence of constants  $A_n \geq 1$  and constant  $c \geq 1$  such that

$$\sum_n A_n \delta_n < \infty, \quad L_n = |y_{n+1} - x_n| + 1 \leq A_1 \dots A_n \cdot c^n. \quad (3)$$

Then  $\sigma(H^\theta) \cap I = \sigma_{pp}(H^\theta) \cap I$  almost everywhere (a.e.) in  $\theta$ .

In the particular case, which is important for many applications, when  $x_n \leq c^n$  (coordinates of blocks do not increase faster than exponential) conditions (2) and (3) become simple. It is enough to require  $\sum_n \delta_n < \infty$ .

Effective verification of (2) and (3) gives the next result. It is based on the physical idea of the absence of resonance.

**Theorem 2.** Suppose (under the conditions of Theorem 1) that for every  $\lambda \in \mathbb{R}^1$  there exist a positive constant  $\varrho = \varrho(\lambda)$  and a sequence of non-resonant blocks  $[x_n, y_n](\lambda)$ , such that  $\text{dist}\{\lambda, \sigma(H^{(n)})\} \geq \varrho(\lambda)$  and  $l_n = |x_n - y_n| + 1 \geq A(\varrho) \ln n$  where the function  $A = A(\varrho)$  could be specified. Then, if  $x_n < c^n$ ,  $c > 0$ ,

$$\sigma(H^\theta) = \sigma_{pp}(H^\theta) \quad \text{a.e. in } \theta.$$

Theorems 1 and 2 include all well-known mechanisms of localization: high barriers, long bumps, gaps in periodic potentials, etc. (See for comparison [9].)

The typical examples where Theorems 1 and 2 could (and will) be applied are random (usually nonhomogeneous) potentials  $V(x, w)$ ,  $x \in \mathbb{Z}^1$ .

**Example 1** (Independent Random Variables). Let  $V(x) = \xi_x(w)$ ,  $x \geq 0, w \in (\Omega, F, P)$  be i.r.v and  $\xi_x(w) = a(x) + \eta_x(w)$ ,  $x > 0$ , where  $a(x)$  is an arbitrary

nonrandom function and  $\{\eta_x(w), x \geq 0\}$  is uniformly non-degenerate in the following sense: There exist positive constants  $\varepsilon_0, \delta_0$  such that

$$\underline{P}\{\eta_x > 1 + \delta_0\} \geq \varepsilon_0, \quad P\{\eta_x < -1 - \delta_0\} \geq \varepsilon_0, \quad x = 0, 1, \dots .$$

Then  $\sigma(H^\theta = A + \xi_x) = \sigma_{pp}(H^\theta)$  a.s.  $\underline{P}$  and a.e. in  $\theta$ .

**Example 2** (Gaussian Potentials). Let  $V(x) = \xi_x$ ,  $x \geq 0$ , be a nonstationary gaussian sequence,  $\langle \xi_x \rangle = 0$ ,  $0 < c_1 < \langle \xi_x^2 \rangle < c_2 < \infty$ ,  $\text{cov}(\xi_x, \xi_y) = \langle \xi_x \xi_y \rangle \leq \frac{c_3}{\ln^{1+\varepsilon}(|x-y|)}$ ,  $|x-y| \geq c_4$ . The last condition is well known in the theory of gaussian fields and processes: if correlations decay only logarithmically, then the structure of the high peaks suffers from bifurcations. It's possible to prove that in our case,  $P$ -a.s. in any interval  $[2^n, 2^{n+1})$ ,  $n > n_0(w)$ , there exists a “triplet”  $|\xi_x| > \delta \sqrt{n}$ ,  $|\xi_{x+1}| > \delta \sqrt{n}$ ,  $|\xi_{x+2}| > \delta \sqrt{n}$  of the high peaks. It's enough to allow us to apply Theorem 2 and so  $\sigma(H^\theta) = \sigma'_{pp}(H^\theta)$  a.e. in  $\theta$  and  $P$ -a.s.

**Example 3** (Unbounded Quasiperiodic Potentials). Let  $f(t) = f(t+1)$  be periodic function on the unit circle  $S^1$  for which there exists at least one point  $t_0$  of logarithmic singularity:

$$|f(t)| \geq c(\ln |t - t_0|^{-1})^{1+\varepsilon}, \quad t \in S^1, \quad \varepsilon > 0. \quad (4)$$

Then for quasirandom potentials of the form

$$V(x) = f(\alpha x), \quad V(x) = f(\alpha x^2 + \beta) \quad (5)$$

for almost all values of parameters  $\alpha$ ,  $(\alpha, \beta)$  we do have a pure point spectrum (a.e. in  $\theta$ ). The second of these potentials (for the case  $f(t) = \text{ctg } \pi t$ ) is the popular model of “quantum chaos” (see [18]).

For the case of entire axis  $\mathbb{Z}^1$ , the treatment for the operator

$$H = A + V(x), \quad -\infty < x < \infty$$

is more subtle. Of course, if conditions of Theorems 1 and 2 are satisfied for  $x > 0$  and  $x < 0$ , then it's possible to prove (as in Theorems 1 and 2) that for a.e.  $\lambda \in I$  (or  $\lambda \in R^1$ )

$$R_\lambda(0, x) = (H - \lambda)^{-1}(0, x) \in l^2(\mathbb{Z}^1).$$

In this spectral problem there is no exterior random parameter (such as  $\theta$  in the case  $\mathbb{Z}_+^1$ ) and (following [5] and [6]) the randomness must be included in the potential. The next lemma generalizes the one-dimensional results of [5] and [6].

**Lemma 1.** Let  $H^a = A + V_0(x) + a\varphi(x)$ ,  $x \in \mathbb{Z}^1$ , for a.e.  $\lambda < R^1$  and denote by

$$R_\lambda^0(0, x) = (A + V_0 - \lambda)^{-1}(0, x) = (H^0 - \lambda)^{-1}(0, x) \in l^2(\mathbb{Z}^1).$$

If perturbation  $\varphi$  decays “sufficiently fast”, then for a.e.  $\lambda$  any two Weyl's solutions  $\psi_\lambda^\pm(x)$  of the equation  $(H^0 - \lambda)\psi = 0$  satisfy

$$\sum_{x \in \mathbb{Z}^1} |\psi^-| |\psi_\lambda^+(x)| |\varphi(x)| < \infty$$

Then, a.e. in  $a$ ,  $\sigma(H^a) = \sigma_{pp}(H^a)$ . The typical  $\psi$  is the one which satisfies  $|\varphi(x)| \leq \exp(-\varepsilon|x|)$ .

In [5] and [6] similar results were proved for functions  $\varphi$  with finite support. We will give two examples of the above lemma.

**Example 4.** Let  $V(x) = \xi_x$ ,  $x \in \mathbb{Z}^1$  be i.r.v. which satisfy the conditions of Example 1. If, in addition, one of them, say  $\xi_0$ , has absolutely continuous distribution, then

$$\sigma(\mathcal{A} + \xi_x) = \sigma_{pp}(\mathcal{A} + \xi_x) \quad P\text{-a.s.}$$

**Example 5.** Let  $\xi_x$ ,  $x \in \mathbb{Z}^1$  be i.i.d.r.v. with common absolutely continuous distribution  $\langle \xi_x \rangle < \infty$  and  $\varphi(x)$ , such that  $|\varphi(x)| \leq \exp(-\varepsilon|x|)$ , be the elementary potential. Consider a homogeneous random potential of alloy type

$$V(x) = \sum_{n=-\infty}^{+\infty} \xi_n \varphi(x - n).$$

It's not very difficult to verify the application of Theorem 2 in this case (see [17], where similar problems were analyzed in a more complicated situation). Application of Lemma 1 to the “partition”

$$V(x) = \sum_{n \neq 0} \xi_n \varphi(x - n) + \xi_0 \varphi(x) = V_0(x) + \xi_0 \varphi(x)$$

shows that the Schrödinger operator  $H = \mathcal{A} + \sum_n \xi_n \varphi(x - n)$  under the above conditions has  $P$ -a.s. p.p. spectrum in  $l^2(\mathbb{Z}^1)$ . Earlier (see [17]) it was known only that  $\sigma(H) = \sigma_{\text{sing}}(H)$   $P$ -a.s.

## § 2. Some Generalizations

The main idea of Theorems 1 and 2 can be extended to more general one-dimensional lattice systems: Jacobi operators of the form  $H\psi(x) = l(x-1)\psi(x-1) + l(x)\psi(x+1)$ , operators where the Laplacian  $\mathcal{A}$  is replaced by nonlocal convolution  $\sum_l l(x-y)\psi(y) = A\psi(x)$  and kernel  $l(x-y)$  decays fast enough off the diagonal and so on.

Let us formulate two theorems of this type. The first result about the “random string” has a clear mechanical meaning.

**Theorem 3.** Let  $H^\theta \psi(x) = l(x-1)\psi(x-1) + l(x)\psi(x+1)$ ,  $l(x) > 0$ ,  $x \geq 0$  be the operator of the lattice string with boundary condition  $\psi(-1) \cos \theta + \psi(0) \sin \theta = 0$ ,  $\theta \in [0, \pi]$ . Assume that for some sequence of blocks  $[x_n, y_n]$ :  $l_n = |x_n - y_n| + 1 > c \ln^{1+\varepsilon} n$ ,  $\varepsilon, c > 0$ ;  $x_n < c_1^n$ ,  $c_1 > 1$  we have  $l(x) \leq \lambda_0$ ,  $x \in [x_n, y_n]$ ,  $n = 1, 2, \dots$ . Then a.e.  $\theta \in [0, \pi)$

$$\sigma(H^\theta) \bigcap (|\lambda| > \lambda_0) = \sigma_{pp}(H^\theta) \bigcap (|\lambda| > \lambda_0).$$

Physically this means that “long, not very elastic inclusions in an elastic medium” lead to localization of short waves. In some cases it is possible to prove that for  $|\lambda| < |\lambda_0|$  there is no p.p. spectrum.

The following is connected with nonlocal Laplacian.

**Theorem 4.** Let  $H = \tilde{\Delta} + V(x)$ , where  $\tilde{\Delta}\psi(x) = \sum_{y \in \mathbb{Z}} l(x-y)\psi(y)$ ,  $|l(z)| \leq \frac{c}{1+|z|^\beta}$ ,  $\beta > 8$  and  $V(x)$  is i.i.d.z.v. with common a.c. distribution. Then P-a.s.

$$\sigma(H) = \sigma_{pp}(H).$$

Other generalizations referred to increasing the dimension.

**Schrödinger Operator in the Strip.** Let  $D = \mathbb{Z}_+^1 + \mathbb{Z}_N$ , where  $\mathbb{Z}_N = (0, 1, \dots, N-1)$ ,  $N \equiv 0$ , is a group of residue mod  $N$  and hamiltonian  $H$  in  $l^2(D)$  has a form

$$\begin{aligned} H^\theta &= \Delta + V(x, z), \quad (x, z) \in D, \quad \Delta\psi(x, z) = \psi(x+1, z) + \psi(x-1, z) \\ &\quad + \psi(x, z+1) + \psi(x, z-1), \quad x \geq 1, \quad z \pm 1 = z \pm 1 \pmod{N}, \\ &\quad \cos \theta_z \psi(-1, z) + \sin \theta_z \psi(0, z) = 0, \quad z \in \mathbb{Z}_N, \\ &\quad \theta = (\theta_0, \dots, \theta_{N-1}) \in [0, \pi]^N. \end{aligned} \tag{6}$$

It is well known, that for homogeneous fields  $V(x, z)$  (with group of shifts  $x \rightarrow x+h$ ,  $x, h \in \mathbb{Z}^1$ ) the localization theorems may be proven by the classical method of Ljapunov exponents. This method in the case of the strip is more complicated than for  $\mathbb{Z}^1$ . On the contrary, the cluster method of Section 1 doesn't feel the difference between  $\mathbb{Z}^1$  and  $\mathbb{Z}^1 \times \mathbb{Z}_N$ .

**Theorem 5.** Assume that for given operator  $H^\theta$  and energy interval  $I \subset R^1$  there exists a system of blocks  $B_n = \{(x, z) : (x_n \leq x \leq y_n, z \in \mathbb{Z}_n), n = 1, 2, \dots$  and constants  $h_n, \varrho_n$ , such that

- a)  $|V(x, z)| \geq h_n$ ,  $(x, z) \in B_n$
- b)  $\text{dist}(h_n, I) = 4 + \varrho_n$ ,  $\varrho_n > 0$ .

If  $x_n < c^n$ ,  $c > 1$ ,  $l_n = |x_n - y_n| + 1$  and  $\sum_n (1 + \varrho_n)^{-l_n} < \infty$ , then

$$\sigma(H^\theta) \cap I = \sigma_{pp}(H^\theta) \cap I \quad \text{a.e. } \theta \in [0, \pi]^N.$$

Note that the central conditions

$$\text{dist}(h_n, I) \geq 4 + \varrho_n, \quad \varrho_n > 0, \quad n = 1, 2, \dots$$

have a slightly different form with respect to “pure one-dimensional” theory. Early on we used  $\text{dist}(h_n, I) > 2 + \varrho_n$ . This is because  $\sigma(\Delta) = [-2, +2]$  in  $l^2(\mathbb{Z}^1)$  and  $\sigma(\Delta) = [-4, 4]$  in  $l^2(\mathbb{Z}^1 \times \mathbb{Z}_N)$ .

In virtue of Theorem 5 and its generalization to the case of  $(\mathbb{Z}^1 \times \mathbb{Z}_N)$ , which uses simple variants of Lemma 1, all previous examples (1-5) automatically transfer to the corresponding examples in the strip (half-strip). (The number of a.c. conditions now equals  $N$ . For example, in the analog of Lemma 1 we must

consider  $N$  perturbations:  $\xi_1\varphi_1(x, z_1), \xi_2\varphi_2(x, z_2), \dots, \xi_N\varphi_N(x, z_N)$ ;  $(\xi_1, \dots, \xi_N)$  has a.c. distribution.)

Multidimensional generalizations ( $d > 1$ ) of Theorems 1 and 2 of Section 1 exist, but are noneffective in the case of homogeneous potentials. They act as follows: For any  $n = 1, 2, \dots$  in  $\mathbb{Z}^d$  let there exist two 1-connected (in the sense of percolation theory) surfaces  $\Gamma_n^+, \Gamma_n^-$ ,  $D_n = \text{diam } \Gamma_n^-, d_n = \min|x - y|_{x \in \Gamma^+, y \in \Gamma^-}$ ,  $\text{Int } \Gamma_1^- \subset \text{Int } \Gamma_1^+ \subset \text{Int } \Gamma_2^- \subset \text{Int } \Gamma_2^+ \subset \dots$ . Let  $|\Gamma_n^\pm| = \underline{\Omega}(D_n^{d-1})$ ,  $|\text{Int } \Gamma_n^\pm| = \underline{\Omega}(D_n^d)$ ,  $n \rightarrow \infty$  and for potential  $V(x)$ ,  $x \in \mathbb{Z}^d$  the next estimation takes place:

$$|V(x)| \geq c(D_n^{d-1} n^{1+\varepsilon})^{1/d_n}, \quad x \in B_n = \text{Int } \Gamma_n^+ \setminus \text{Int } \Gamma_n^-.$$

Under these conditions the cluster expansion of resolvent  $R_\lambda^0(0, x) = (\Delta + V(x))^{-1}(0, x)$ , with respect to non-resonant families of blocks  $B_n(\partial B_n = \Gamma_n^+, \Gamma_n^-)$ , shows that

$$R_\lambda(0, x) \in l^2(\mathbb{Z}^d) \quad \text{a.e. } \lambda.$$

Using the Simon-Wolff theorem under some additional technical conditions it is possible to prove a few concrete results about the p.p. spectrum.

**Example 6.** Let  $V(x) = \xi_x|x|^\alpha$ ,  $x \in \mathbb{Z}^d$ ,  $\alpha > 0$  and  $\xi_x$  is i.i.d.r.v. with some moments properties and a.c. distribution. For example,  $\xi_x \in [-1, 1]$  or  $[0, 1]$  and is uniformly distributed. Another case is where  $\xi_x$  is standard  $N(0, 1)$  gaussian r.v. Here  $P$ -a.s. for all  $\alpha > 0$

$$\sigma(\Delta + V(x)) = \sigma_{pp}(\Delta + V(x)).$$

The corresponding eigenfunctions decay superexponentially. The structure of the  $\sigma_{\text{ess}}$  in the case  $\xi_x \in [0, 1]$  is not trivial and dependent on  $\alpha$ . Note that  $\sigma(\Delta + V(x)) = \sigma_{\text{discr}}(\Delta + V(x))$  if  $\alpha > d$ .

**Example 7.** There exist potentials  $V(x, w)$ ,  $x \in \mathbb{Z}^d$ ,  $d > 1$  such that

a)  $V(x, w)$  is homogeneous and ergodic with “good” mixing properties. For example it has strong mixing condition with respect to the family of bounded subsets  $\mathbb{Z}^d$ .

b)  $\langle |V(x, w)|^p \rangle < \infty$  for every  $p > 0$ .

c) Operator  $H = \Delta + \sigma V(x)$  has p.p. spectrum for all positive coupling constants  $\sigma$ .

Note, however, that in this “counterexample” to the Anderson’s hypothesis the correlations of  $V(x, w)$  decay (in any sense) very slowly. Potentials  $V(x, w)$  of this example do not percolate from above. That is, for any  $\lambda > 0$  the set  $A^-(\lambda) = \{x : |V(x, w)| < \lambda\}$  is the union of finite components. Typical potentials (in physical applications) have a finite level of percolation [16].

### § 3. Continuous One-Dimensional Case

In the transition from the results of Section 1 for lattice Laplacians to their analogs for the operator

$$H^\theta \psi(x) = -\frac{d^2\psi}{dx^2} + V(x), \quad x \geq 0; \quad \psi(0) \cos \theta + \psi'(0) \sin \theta = 0$$

there are a few technical obstacles connected with unboundness of  $-\Delta = -d^2/dx^2$  in  $l^2(R_+^1)$ . If  $V(x) \geq 0$  the methods and results are similar to Theorems 1–3 but seem stronger.

**Theorem 6.** *Let potential  $V(x) \geq 0$  for some  $\delta_0 > 0$  have a next estimation from below. For some sequence  $x_n \uparrow \infty$ ,  $x_{n+1} - x_n \uparrow \infty$*

$$V(x) > h_n, \quad x \in [x_n, x_n + \delta_0]. \quad (8)$$

If

$$\overline{\lim}_{n \rightarrow \infty} \frac{x_{n+1}}{x_n} \exp(-\delta_0 \sqrt{h_n}) < 1, \quad \sum_n \exp(-\delta_0 \sqrt{h_n}) < \infty \quad (8')$$

then

$$\sigma(H^\theta) = \sigma_{pp}(H^\theta) \quad \text{a.e. } \theta \in [0, \pi).$$

A. Gordon has analyzed in detail the particular case  $V(x) \equiv 0$ ,  $x \notin \cup_n [x_n, x_n + \delta_0]$ ,  $V(x) = h_n \uparrow \infty$ ,  $x \in [x_n, x_n + \delta_n]$  (near high scatterers). He proved that in this case under condition

$$\underline{\lim}_{n \rightarrow \infty} \frac{x_{n+1}}{h_n \cdot x_n} \exp(-\delta_0 \sqrt{h_n}) > 1 \quad (9)$$

$$\text{Sp } H^\theta = \text{Sp}_{sc} H^\theta \quad \text{for a.e. in } \theta.$$

The next result is a variant of Kotani's theorem [3], but for p.p. spectrum (but not singular as in [3]).

**Theorem 7.** *Let  $F_0(x) = F_0(x + T)$  be a continuous periodic function and  $[x_n, y_n]$  be a sequence of blocks such that  $\sup_{x \in [x_n, y_n]} |V(x) - F_0(x)| = \varepsilon_n \rightarrow 0$   $\underline{\lim}_{n \rightarrow \infty}$ ,  $x_n < c^n$ ,  $c > 1$ ,  $l_n = |x_n - y_n| \geq \ln^{1+c} n$ ,  $\varepsilon > 0$ ,  $n \geq n_0$ . If  $\Delta$  is one of the gaps in the  $\sigma_{ess}(-d^2/dx^2 + F_0(x))$ , then for a.e. in  $\theta$*

$$\sigma(H^\theta) \cap \Delta = \sigma_{pp}(H^\theta) \cap \Delta.$$

If potential  $V(x)$  is unbounded from below, but has a logarithmic estimation of the form  $V(x) \geq -c_1 \ln(|x| + 1) + c_2$ , then it is possible to prove variants of Theorems 5 and 6 under stronger assumptions on the increase of  $\{h_n\}$  in Theorem 5 or  $\{l_n\}$  in Theorem 6.

In the case of the entire line  $R^1$  (that is in  $l^2(R^1)$ ) the continuous analog of Lemma 1 (Sect. 1) can be used. It makes it possible to prove the theorems on the point spectrum for many physically interesting models.

**Example 8.** Let  $V(x) = \sum_{n=-\infty}^{+\infty} \xi_n \varphi\left(\frac{x-x_n}{\theta_n}\right)$  be a “shot noise” potential. Here  $\{x_n\}$  is the Poissonian points flow,  $|\varphi(x)| \leq \exp(-\varepsilon|x|)$  the elementary potential,  $\{\theta_n, \xi_n\}$  i.i.d.r. vectors with finite exponential moments:  $\langle \exp(z\xi) \rangle < \infty$ ,  $\langle \exp(\theta z) \rangle < \infty$ ,  $\langle \exp(z \frac{1}{\theta}) \rangle < \infty$ ,  $|z| < z_0$ ,  $z_0 > 0$  and the distribution of  $\xi_n$  is a.c. Then  $P$ -a.s. in  $l^2(R^1)$

$$\sigma(H) = \sigma(-d^2/dx^2 + V(x, w)) = \sigma_{pp}(H).$$

This example is closely connected with the paper [17] in which the authors proved that in the same situation  $\sigma(H) = \sigma_{\text{sing}}(H)$   $P$ -a.s., but under additional restrictions, the elementary potential  $\varphi$  is not the soliton. Although this condition is not essential and the spectrum is p.p., the appearance of the soliton in this context is not accidental.

**Theorem 8.** *There exists potential  $V(x, w)$ ,  $x \in R^1$  and large parameter  $L > 0$  for which*

$$\text{a)} \quad \left| V(x, w) - \sum_{A=-\infty}^{+\infty} -\left( \frac{2\xi_n^2}{\text{ch}^2 \xi_n (x - l_n)} \right) \right| \leq e^{-\delta L}, \quad \delta > 0, \quad x \in R^1. \quad (10)$$

Here  $\{\xi_n\}$  is i.i.d.r.v. (for instance uniformly distributed in  $[0, a]$ ,  $a > 0$ ),  $\{l_n\}$  is a homogeneous random point process in  $R^1$  (dependent on  $\{\xi_n\}$ ) with good mixing properties. Note that  $\varphi(x) = -2/\text{ch}^2 x$  is the simplest soliton (1-soliton).

$$\text{b)} \quad \sigma(H) \cap [0, \infty) = \sigma_{\text{ac}}(H) \cap [0, \infty), \quad \sigma(H) \cap [-\infty, 0] = \sigma_{pp}(H) \cap [-\infty, 0].$$

This potential is one of the realizations of “soliton’s gase”. It is closely connected with the problem of statistical solutions of the KdV-equation.

## § 4. Parabolic Problems for the Anderson Model. Intermittency and Related Topics

Evolution problems for the physical fields in the random medium (chemistry kinetics, hydrodynamics, etc.) very often have the form of a parabolic equation with random coefficients, in particular, with random potential. The simplest example is

$$\begin{aligned} \partial c / \partial t &= D \Delta c + \xi(x) c \\ c(0, x) &\equiv 1. \end{aligned} \quad (11)$$

Function  $c(t, x)$ ,  $t \geq 0$  has the meaning of the concentration of the particles at moment  $t$  in the point  $x$ . The kinematic part of the hamiltonian  $D\Delta$  describes its diffusion ( $D$ -diffusion coefficient) and potential  $\xi(x)$  its transformation. If  $\xi(x) > 0$ , then  $\xi(x)dt$  is the probability that in the time interval  $(t, t+dt)$  any particle in the point  $x$  will split (birth of a particle). If  $\xi(x) < 0$ , then  $\xi(x)$  is the intensity of the death process in the point  $x$ . We will consider problem (11) in the discrete case where  $x \in \mathbb{Z}^d$ ,  $D\Delta\psi(x) = D \sum_{|x'-x|=1} (\psi(x') - \psi(x))$  is the lattice

Laplacian (which is the generator of the symmetrical random walk  $x_t, t \geq 0$  in continuous time) and  $\xi(x, w)$  is the homogeneous ergodic random field.

If  $V(x) = \xi_x$ ,  $x \in \mathbb{Z}^d$  are i.i.d.r.v., then  $H - D\Delta + \xi$  is the hamiltonian of the tight-binding Anderson model. The diffusion coefficient  $D$  is the inverse coupling constant  $D = 1/\sigma$ . It's well known [7], [8], that for  $\sigma \gg 1$  (strong disorder) or for small  $\sigma$ , but  $\lambda \gg 1$  (fluctuation part of the spectrum) under some technical restrictions  $P$ -a.s.  $\sigma(H) = \sigma_{pp}(H)$ .

**Asymptotic Properties.** The solution  $c = c(t, x)$ ,  $t \rightarrow \infty$  can be represented in the spectral terms. It allows investigation by direct probabilistic methods. This gives us essential information about the structure  $\sigma(H)$ ,  $\lambda \gg 1$ .

The central qualitative property of the field  $c(t, x)$ ,  $t \rightarrow \infty$  is its intermittency – that is, informally, the existence of strongly pronounced spatial structures (in this case sharp and high peaks). The definition of intermittency is given in terms of statistical moments  $c(t, x)$ .

Remember that  $c(0, x) \equiv 1$ . This means that  $c(t, x)$  is a homogeneous ergodic field for every  $t > 0$ . It is not very difficult to prove that the condition  $\langle \exp(t\xi) \rangle < \infty$ ,  $t > 0$  is necessary and sufficient for the existence of  $\langle c^p(t, x) \rangle$ ,  $t > 0$ ,  $p = 1, 2, \dots$ .

**Definition 1.** We will say that the family of the fields  $c(t, x)$ ,  $x \in \mathbb{Z}^d$ ,  $t > 0$  is the asymptotic intermittency parameter of the family when  $t \rightarrow \infty$ , if for the functions

$$\Lambda_p(t) = \ln \langle c^p(t, \cdot) \rangle$$

the following relations take place as  $t \rightarrow \infty$

$$\Lambda_t(t) \ll \frac{\Lambda_2(t)}{2} \ll \frac{\Lambda_B}{B}(t) \ll \dots$$

Here  $A(t) \ll B(t)$  means  $B(t) - A(t) \rightarrow_{t \rightarrow \infty} +\infty$ .

**Theorem 9.** If potential  $\xi(x)$  is unbounded from above and  $G(t) = \langle \exp(t\xi(\cdot)) \rangle < \infty$ ,  $t > 0$ , then the solution  $c(t, x)$  is asymptotically intermittent in the sense of Definition 1. The logarithmic asymptotics of the statistical moment has a form

$$\frac{\ln \langle c^p(t, 0) \rangle}{p} = \frac{\Lambda_p(t)}{p} \sim_{t \rightarrow \infty} \frac{G(pt)}{p}. \quad (12)$$

The more exact asymptotics depends upon the structure of the tails of the one-dimensional distributions of potential  $\xi(\cdot)$ .

**Theorem 10.** Let  $\xi(x)$ ,  $x \in \mathbb{Z}^d$  be i.i.d.r.v. (Anderson model) and  $P\{\xi(1) > t\} \sim_{t \rightarrow \infty} \exp(-ct^\beta)$ ,  $\beta > 1$ . Then

$$\frac{\Lambda_p(t)}{p} = \frac{G(pt)}{p} - 2d Dpt + \bar{\bar{\Omega}}(t) = \frac{c(\beta)t^{\beta/\beta-1}}{p} - 2d Dpt + \bar{\bar{\Omega}}(t), \quad t \rightarrow \infty. \quad (13)$$

But if  $P\{\xi(\cdot) > t\} \sim \exp(-c \exp(c, t^\beta))$ ,  $\beta > 1$ , then

$$\frac{\Lambda_p(t)}{p} = \frac{G(pt)}{p} + \bar{\bar{\Omega}}(t) = c(\beta, p) \ln^{\frac{1}{\beta}} t + \bar{\bar{\Omega}}(t). \quad (14)$$

The difference between these two formulas are due to physical reasons. In the first case “strong centers” of the potential  $\zeta(\cdot)$ , which contribute mainly to the growth of the number of particles, have the form of a single high peak. In the second case it is wide but not very high islands. The second term of asymptotics of  $\wedge_p(t)/p$  describes the probability of “keeping” particles by “strong centers”.

It is possible to observe the same effect in the results about the almost surely (a.s.) behavior of  $c(t, x)$ ,  $t \rightarrow \infty$ .

**Theorem 11.** Let  $\ln \ln(P\{\xi > t\})^{-1} < c_1 t^\beta$ ,  $\beta < 1, c_1 > 0$  (roughly speaking,  $P\{\xi > t\} \geq \exp\{-c \exp(c_1 t^\beta)\}$ ,  $\beta < 1$ ). Then  $P$ -a.s.  $t \rightarrow \infty$

$$\frac{\ln c(t, x)}{t} = \left( \log \left( \frac{1}{H(0)} \right) \right)^{-1} (d \ln t) - 2dD + \bar{O}(1). \quad (15)$$

(Here,  $H(t) = P\{\xi > t\}$ ,  $(\ )^{-1}$  means inverse function.)

If, however,

$$\ln \ln H(t) > c_2 t^\beta, \quad \beta > 1, \quad c_2 > 0 \quad (16)$$

then  $P$ -a.s.

$$\frac{\ln c(t, x)}{t} = \left( \log \frac{1}{H(\cdot)} \right)^{-1} (d \ln t). \quad (17)$$

Consider now the initial parabolic Anderson problem (11) for the localized initial condition  $c(0, x) = \delta_0(x)$ . Assume, that  $\xi(x) > 0$  is an i.i.d.r.v. with “exponential tails”. For simplicity let  $P\{\xi > t\} = \exp(-ct^\beta)$ ,  $\beta > 1$ . In the beginning we have only one particle but birth and diffusion processes lead to the “occupation of the space”. This problem of the quantitative description of this phenomena is similar to the famous problem KPP (Kolmogorov-Petrovski-Piskunov).

It is not very difficult to show that the boundary of “occupied” region in moment  $t$  is given by a sphere

$$S_t = \{x : |x| \leq t \ln^{1/\beta} t\}$$

in the following way: if  $|x| > t \ln^{1/\beta+\varepsilon} t$ , then  $c(t, x) \xrightarrow{P} 0$  and if  $|x| < t \ln^{1/\beta-\varepsilon} t$ , then  $c(t, x) \xrightarrow{P} \infty$ .

However, the set  $\sigma_t$  is extremely nonuniformly occupied, as is a typical effect of intermittency.

**Theorem 12.** For every  $t > 0$ ,  $\varepsilon > 0$  it's possible to find (random) points  $x_1(t, w), \dots, x_k(t, w)$ ,  $k < \ln t$ , such that

$$\sum_{x \in \mathbb{Z}^d} c(t, x) \geq (1 - \varepsilon) \sum_{i=1}^k c(t, x_i). \quad (18)$$

It is likely that  $k = k(t, w)$  is bounded in probability. This theorem shows that it is necessary to be very careful in the applications of the results on the averaging description of the propagation front of concentration in random media. The field of concentration inside the front has an extremely non-uniform structure.

## § 5. On the Basic States in the Anderson Model. Precision of the Asymptotical Formulas for “Lifshitz Tails”

The limit theorems for the boundary part of the spectrum, that is for the basic states will be considered in finite, but big volume  $V$  when  $V \rightarrow \infty$ . The integral density of states  $N(\lambda)$  can be studied in the framework of the same procedure.

Let  $S_N^d = [-N, N]^d$ , where points  $N, -N$  are identified, be the  $d$  dimensional lattice torus of the volume  $V_N = (2N)^d$  and  $H = \Delta + \xi(x)$  be the operator of the tight-binding Anderson model. This means that  $\xi(x)$ ,  $x \in S_N^d$  is an i.i.d.r.v. We will consider a typical one for the theory of “Lifshitz tails” case

$$P\{\xi_x > t\} = \exp\{-ct^\beta\}, \quad t > 0, \quad \beta > 0.$$

We are interested in the structure of higher eigenvalues and corresponding eigenfunctions. It's easy to understand that they are closely connected with the higher peaks of the potential  $\xi(x)$ ,  $x \in S_N^d$ . Consider the two variational series:

$$\begin{aligned} \xi_N^{(1)} &> \xi_N^{(2)} > \dots > \xi_N^{(V)} \\ \lambda_N^{(1)} &> \lambda_N^{(2)} > \dots > \lambda_N^{(V)} \end{aligned}$$

The limit distribution of any fixed number  $k$  of the first r.v. in the  $\xi$ -series is described by the Weibull's type law:

$$\begin{aligned} P\left\{ \frac{\xi_N^{(1)} - \xi_N^{(2)}}{\frac{1}{p} \ln^{(1/p-1)} V_N} \in (x_1 + dx_1), \frac{\xi_N^{(2)} - \xi_N^{(3)}}{\frac{1}{p} \ln^{(1/p-1)} V_N} \in (x_2 + dx_2), \dots \right. \\ \left. \frac{\xi_N^{(k)} - \ln^{1/p} V_N}{\frac{1}{p} \ln^{(1/p-1)} V_N} \in (x_k + dx_k) \right\} \\ \rightarrow_{N \rightarrow \infty} P_k(x_1, \dots, x_k) = \exp(-x_1 - 2x_2 - \dots - kx_k - e^{-x_k}), \quad x_1, \dots, x_k \geq 0. \end{aligned} \quad (19)$$

It seems very reasonable, that the corresponding (or near) formulas are valid for  $\lambda_N^{(i)}$ ,  $i = 1, 2, \dots, k$ . In some sense it is true.

**Theorem 13.** For some positive  $A_N^{(1)}, \dots, A_N^{(k)}, B_N$ ,

$$\begin{aligned} P\left\{ \frac{\lambda_N^{(1)} - \lambda_N^{(2)} - A_N^{(1)}}{B_N} > x_1, \frac{\lambda_N^{(2)} - \lambda_N^{(3)} - A_N^{(2)}}{B_N} > x_2, \dots, \frac{\lambda_N^{(k)} - A_N^{(k)}}{B_N} > x_k \right\} \\ \rightarrow_{N \rightarrow \infty} \int_{x_1}^{\infty} \dots \int_{x_k}^{\infty} \exp(-y_1 - 2y_2 - \dots - ky_k - \exp(-y_k)) dy_1 \dots dy_k, \\ x_1, \dots, x_k > 0. \end{aligned} \quad (20)$$

The structure of normalizing constants, however, depends on  $\beta$ . There exist many bifurcations with respect to  $\beta$ . The nature of these bifurcations is very simple. It is obvious that only  $\xi$ -peaks which have an order  $\underline{O}(\ln^{1/\beta} V)$ , for example, bigger than  $(1 - \varepsilon) \ln^{1/\beta} V$  can contribute to the initial part of the  $\lambda$ -series. But  $\#\{x \in S_N^d : \xi(x) > (1 - \varepsilon) \ln^{1/\beta} V\} = \underline{O}(V^{\delta(\varepsilon)})$ ;  $\delta(\varepsilon) \rightarrow 0$ ,  $\varepsilon \rightarrow 0$  and the

distances between these  $\xi$ -peaks have an order  $N^{1-\delta_1(\varepsilon)}$ ;  $\delta_1(\varepsilon) \rightarrow 0$ ,  $\varepsilon \rightarrow 0$ . The interaction between peaks is very small and any of them as shown by standard perturbation calculations gives an eigenvalue

$$\lambda(x_0) = \xi(x_0) - \frac{c_0}{\xi(x_0)} + c_1 \frac{\sum_{|x'-x|=1} \xi(x')}{\xi^2(x_0)} + \bar{O}\left(\frac{1}{\xi^3(x_0)}\right). \quad (21)$$

If  $p < 2$ , then the second and all other terms are small enough with respect to the “gaps” between neighboring  $\xi^{(i)}$  and we can use formula (19), changing  $\xi^{(i)}$  by  $\lambda^{(i)}$ . If  $2 \leq p < 3$ , then the second term of expansion (21) is essential. It is necessary to slightly change constant  $A_N^{(k)}$ . However, as in the case  $p < 2$ , we have the correspondence  $\xi^{(1)} \leftrightarrow \lambda^{(1)}$ ,  $\xi^{(2)} \leftrightarrow \lambda^{(2)}$ , ...,  $\xi^{(k)} \leftrightarrow \lambda^{(k)}$ . If  $p \geq 3$ , then the gaps between  $\xi^{(1)} - \xi^{(2)}, \dots$  are smaller than  $\frac{1}{\xi^2(x_0)}$ , this correspondence is destroyed and all normalized constants are new. It is not possible to write out explicit formulas for  $A_N^{(i)}, B_N$  as functions of  $\beta$  and dimension  $d$ .

The same analysis applies to the problem of “Lifshitz tails” as to the high energy asymptotic of  $N^*(\lambda) = 1 - N(\lambda) = \lim_{V \rightarrow \infty} \frac{1}{V} \sum_{\lambda_N^{(0)} > \lambda} 1$ . It is well known that

$$P\{\xi > \lambda + 2d\} < N^*(\lambda) < P\{\xi > \lambda - 2d\}.$$

(Of course,  $2d = \|A\|_p$ ). It follows from these estimations that

$$-\ln N^*(\lambda) \sim_{\lambda \rightarrow +\infty} \lambda^\beta.$$

What is the more precise form of this asymptotics? The answer depends primarily on  $\beta$  and contains information on the structure of corresponding eigenfunctions.

#### Theorem 14.

a) If  $p < 2$ , then

$$-\ln N^*(\lambda) = \lambda^\beta + \bar{O}(1).$$

b) If  $p = 2$  (gaussian case), then

$$-\ln N^*(\lambda) = \lambda^\beta + c_1(d) + \bar{O}(1). \quad (22)$$

c) If  $2 < p < 3$

$$-\ln N^*(\lambda) = \lambda^\beta + c_2 \lambda^{\beta-2} + \bar{O}(1). \quad (23)$$

d) If  $3 < p < 4$

$$-\ln N^*(\lambda) = \lambda^\beta + c_2 \lambda^{\beta-2} + c_3 \lambda^{(\beta-2)-\frac{2}{p-1}} (1 + \bar{O}(1)) \quad (24)$$

and so on.

There is no room here to discuss other types of bifurcations of  $N(\lambda)$  in the case of “double exponential” tails of the distribution  $\xi(\cdot)$ . The situation here is similar to the results of Theorems 10 and 11.

## References

- [1] Kotani, S.: (1983), Ljapunov indices determine absolute continuous spectra of stationary one-dimensional Schrödinger operators. In: Proc. Taniguchi Intern. Symp. on Stochastic Analysis, Katata and Kyoto (1982), ed. K. Ito. North-Holland, pp. 225–247
- [2] Kotani, S. (1986): Ljapunov exponent and spectra for one dimensional random Schrödinger operators, Proc. Conf. on Random Matrices and their Applications. Contemp. Math. **50**, Providence R.I., pp. 277–286
- [3] Kotani, S. (1985): Support theorems for random Schrödinger operators. Comm. Math. Phys. **97**, 443–452
- [4] Fröhlich, G., Spencer, T. (1983): Absence of diffusion in the Anderson tight-binding model for large disorder or low energy. Comm. Math. Phys. **88**, 151–189
- [5] Delyon, F., Levy, Y., Souillard, B. (1985): Anderson localization for multidimensional systems at large disorder or low energy. Comm. Math. Phys. **100**, 4670–470
- [6] Simon, B., Wolff, T. (1986): Singular continuous spectrum under rank one perturbations and localization for random Hamiltonian. Commun. Pure Appl. Math. **39**, 75–90
- [7] Martinelli, F., Scoppola, E. (1987): Introduction to the mathematical theory of Anderson localization. La Rivista del Nuovo Cimento **10**, 3
- [8] Dreyfus, H. von, Klein, A. (1989): A new proof of localization in the Anderson tight-binding model. Comm. Math. Phys. **124**, 285–299
- [9] Simon, B., Spencer, T. (1989): Trace class perturbations and the absence of absolutely continuous spectra. Comm. Math. Phys. **125**, 113–125
- [10] Kirsch, W., Molchanov, S., Pastur, L. (1989): The point spectrum of the one-dimensional Schrödinger operator with an unbounded potential. Preprint Ruhr-Universität, Bochum
- [11] Kirsch, W., Molchanov, S., Pastur, L. (1990): One-dimensional Schrödinger operator with unbounded potential: pure point spectrum, I. Funct. Anal. Appl. **2** (in Russian)
- [12] Kirsch, W., Molchanov, S., Pastur, L. (1990): One-dimensional Schrödinger operator with unbounded potential: continuous case, II. Funct. Anal. Appl. (to appear) (in Russian)
- [13] Kirsch, W., Molchanov, S. (1989): Schrödinger operator with the potential of “solitons gas” type. Preprint Ruhr-Universität, Bochum
- [14] Gordon, A., Molchanov, S. (1990): One-dimensional Schrödinger operators with a strongly fluctuating random potential. Funct. Anal. Appl. (to appear) (in Russian)
- [15] Görtner, J., Molchanov, S. (1990): Parabolic problems for the Anderson model. Intermittency and related topics, I. Comm. Math. Phys. (to appear)
- [16] Men'shikov, M., Molchanov, S., Sidorenko, A. (1986): Percolation theory and some applications. Itogi Nauki, Ser. Teor. Veroyatn **25**, 53–110 (in Russian)
- [17] Kirsch, W., Kotani, S., Simon, B. (1985): Absence of absolutely continuous spectrum for one-dimensional random but deterministic Schrödinger operators. Ann. Inst. H. Poincaré **42**, 383
- [18] Casati, G., Chirikov, B., et al. (1979): Stochastic behavior in classical and quantum Hamiltonian systems. Lecture Notes in Physics, vol. 93. Springer, Berlin Heidelberg New York
- [19] Molchanov, S., Simon, B.: Localization theorem for non-local Schrödinger operators in one-dimensional lattice case. Comm. Math. Phys. (to appear)
- [20] Jacksic, V., Molchanov, S., Simon, B.: Spectral theory of multidimensional lattice Schrödinger operator with random potential, which increase in probability. Comm. Math. Phys. (to appear)



# The Laws of Some Brownian Functionals

Marc Yor

Laboratoire de Probabilités, Université P. et M. Curie, Tour 56, 3ème étage  
4, place Jussieu, F-75252 Paris Cedex 05, France

Thanks mainly to the relationship between the heat equation, newtonian potential theory and Brownian motion, the laws of a large number of Brownian functionals have been obtained during the last fifty years, at least via explicit expressions of their Laplace and Fourier transforms. Much pioneering work in this area was done by Paul Lévy.

Gradually, with the development of Itô's stochastic calculus, excursion theory, path decompositions and the technique of enlargement of filtrations, these studies of individual distributions on  $\mathbb{R}$ , sometimes exhibiting identities between two laws, which looked a priori to be mere "coincidences", have been understood in a deeper way, in fact often by showing that two *processes* are identical in law; see Biane [3], for a recent survey in that spirit.

The most elementary examples of Brownian functionals are linear functionals: if  $f \in L^2(\mathbb{R}_+, dt)$ , and  $(B_t, t \geq 0)$  is a real-valued BM, the Wiener integral  $\int_0^\infty f(t) dB_t$  is a centered Gaussian variable, with variance  $\int_0^\infty f^2(t) dt$ . Quadratic functionals of BM represent the next level of complexity; those functionals are of great interest as, somewhat surprisingly, they occur in a number of very different studies of Brownian motion, such as the Ray-Knight theorems for Brownian local times, the Ciesielski-Taylor identities, some limiting laws of planar BM, and principal values of Brownian local times.

We shall take here, as a prototype of a quadratic Brownian functional, the stochastic area of planar BM, and it will be shown how Paul Lévy's formula for this stochastic area appears again and again in most of the above mentioned studies of Brownian motion.

## 1. On Lévy's Area Formula

Consider a two-dimensional Brownian motion  $Z_t = X_t + iY_t$ ,  $t \geq 0$ , starting from 0, and the stochastic area process

$$S_t = \int_0^t (X_s dY_s - Y_s dX_s), \quad t \geq 0.$$

Lévy's formula:

$$E[\exp(i\lambda S_1)|Z_1 = z] = \left( \frac{\lambda}{\sinh \lambda} \right) \exp - \frac{|z|^2}{2} (\lambda \coth \lambda - 1) \quad (1)$$

has played some important rôle in recent years, for example in the Bismut approach [7, 8] to the Atiyah-Singer theorems.

To prove formula (1), Lévy [14] used a diagonalization procedure. A different approach (Williams [25], Yor [28]) is to use a change of probability method, which reduces the computation of the law of the quadratic functional  $S_1$  to that of the variance of a Gaussian variable.

First, by the rotational invariance of the law of BM, and independence properties, we have, for any  $\lambda \in \mathbb{R}$ :

$$E[\exp(i\lambda S_1)|Z_1 = z] = E \left[ \exp - \frac{\lambda^2}{2} \int_0^1 ds |Z_s|^2 |Z_1| = |z| \right]. \quad (2)$$

Next, we introduce the new probability  $P^\lambda$ :

$$P^\lambda|_{\mathcal{F}_t} = \exp \left\{ \frac{\lambda}{2} (|Z_t|^2 - 2t) - \frac{\lambda^2}{2} \int_0^t ds |Z_s|^2 \right\} \cdot P|_{\mathcal{F}_t} \quad (3)$$

under which  $(Z_t, t \leq 1)$  is a Gaussian process (more precisely: an Ornstein-Uhlenbeck process) for which the variance at time 1 is easily computed. Formula (1) now follows from formulae (2) and (3).

We note a simple consequence of formulae (1) and (2):

$$E \left[ \exp - \left( a |Z_1|^2 + \frac{\lambda^2}{2} \int_0^1 ds |Z_s|^2 \right) \right] = \left( \cosh \lambda + 2a \frac{\sinh \lambda}{\lambda} \right)^{-1}. \quad (4)$$

Many variants of Lévy's formula (1) have now been developed; in particular, Biane-Yor [5] obtained a sequence of extensions of Lévy's formula by decomposing  $(Z_u, u \leq 1)$  into the sum of the Brownian bridge  $(Z_u - uZ_1, u \leq 1)$  and  $(uZ_1, u \leq 1)$ , then developing the left-hand side of (1) with respect to this decomposition, and finally iterating this procedure. The Lévy formulae obtained in this way are closely connected, on the one hand, with the linear decomposition of BM along the orthogonal basis of the Legendre polynomials and, on the other hand, with the continued fraction:

$$\lambda \coth \lambda - 1 = \frac{\lambda^2}{3+} \frac{\lambda^2}{5+} \frac{\lambda^2}{7+} \dots$$

## 2. Squares of Bessel Processes and Ray-Knight Theorems

Let  $\delta \geq 1$  be an integer, and  $(q_t, t \geq 0)$  be the square of a BES( $\delta$ ) process, that is the square of the euclidean norm of a  $\delta$ -dimensional BM  $(B_t, t \geq 0)$ ; then,  $q$  satisfies the

SDE:

$$q_t = q_0 + 2 \int_0^t \sqrt{q_s} d\beta_s + \delta t, \quad (5)$$

where  $(\beta_s, s \geq 0)$  is a real-valued BM.

It is well-known that this equation, with  $q_0 = x \geq 0$ , and  $\delta$  any positive real number, has a unique pathwise solution in  $\mathbb{R}_+$ , hence a unique law on  $C(\mathbb{R}_+, \mathbb{R}_+)$ , which we shall denote by  $Q_x^\delta$ .

Shiga and Watanabe [23] remarked that

$$Q_x^\delta * Q_{x'}^{\delta'} = Q_{x+x'}^{\delta+\delta'}, \text{ for any } \delta, \delta', x, x' \geq 0 \quad (6)$$

(where  $P * Q$  is the convolution of  $P$  and  $Q$ ), thus extending to all starting points and dimensions the obvious additivity property for integer dimensions.

From (6),  $Q_x^\delta$  is an infinitely divisible probability distribution on  $C(\mathbb{R}_+, \mathbb{R}_+)$ , which admits the following Lévy-Khintchine representation (Pitman-Yor [18]): there exist two  $\sigma$ -finite  $\geq 0$  measures  $M$  and  $N$  on  $C(\mathbb{R}_+, \mathbb{R}_+)$  such that:

$$Q_x^\delta(\exp - \langle \omega, f \rangle) = \exp - (xM + \delta N)(1 - \exp - \langle \omega, f \rangle) \quad (7)$$

(we use the notation:  $\langle \omega, f \rangle = \int_0^\infty dt \omega(t)f(t)$ , for  $f \geq 0$ ).

This representation may be obtained, including an explicit description of  $M$  and  $N$  in terms of the Itô characteristic measure  $n(d\omega)$  of the Poisson point process of Brownian excursions, with the help of the Ray-Knight theorems on Brownian local times ( $L_t^a$ ;  $a \in \mathbb{R}$ ,  $t \geq 0$ ), which are now presented:

- (RK<sub>1</sub>) if  $T_1 \equiv \inf\{t : B_t = 1\}$ , the law of  $(L_{T_1}^{1-a}; 0 \leq a \leq 1)$  is  $Q_0^2$
- (RK<sub>2</sub>) if  $\tau_x \equiv \inf\{t : L_t^0 = x\}$ , the law of  $(L_{\tau_x}^b; b \geq 0)$  is  $Q_x^0$

(after the original proofs of Ray [22] and Knight [12], several different derivations of (RK<sub>1</sub>) and (RK<sub>2</sub>) have been given; see, for example, Jeulin [11] for a recent survey based on Tanaka's formula).

Conversely, we may now deduce from the Lévy-Khintchine representation (7) some extensions of the Ray-Knight theorems; here is one (Le Gall-Yor [13]): for  $\delta > 0$ , the law of the process  $(C_a^\delta; a \geq 0)$  of the local times of  $(|B_t| + \frac{2}{\delta} L_t^0; t \geq 0)$  is  $Q_0^\delta$ .

In the particular case  $\delta = 2$ , we recover (RK<sub>1</sub>), using the representation of BES(3) as  $(|B_t| + L_t^0; t \geq 0)$ , due to Pitman ([17]), jointly with a well-known time-reversal relation between BES(3) and Brownian motion (see Williams [26]).

### 3. An Explanation of the Ciesielski-Taylor Identities

Ciesielski-Taylor [9] obtained the following puzzling identity in law:

$$\int_0^\infty ds 1_{(R_{\delta+2}(s) \leq 1)} \stackrel{(\text{law})}{=} T_1(R_\delta), \quad (8)$$

where  $R_\delta$  resp.  $R_{\delta+2}$ , is a BES process, with dimension  $\delta$ , resp.  $\delta + 2$ , starting from 0, and  $T_1(R_\delta) \equiv \inf\{t : R_\delta(t) = 1\}$ .

In the case  $\delta = 1$ , the identity (8) can be understood by time-reversal, but no such pathwise explanation has been obtained for other dimensions. Below, a spectral-type explanation and extensions are presented, following [29].

Writing both sides of (8) as integrals of the local times of the two BES processes, and using the Ray-Knight theorems, (8) then appears as a particular case of the identity in law:

$$-\int_0^1 da f'(a) B_{g(a)}^2 \stackrel{\text{(law)}}{=} \int_0^1 da g'(a) B_{f(a)}^2 \quad (9)$$

where  $f, g : [0, 1] \rightarrow \mathbb{R}_+$ , are  $C^1$ , with  $f$  decreasing,  $g$  increasing, and  $f(1) = g(0) = 0$ .

In fact, a more general identity in law holds:

$$-\int_a^b dx f'(x) B_{g(x)}^2 + f(b) B_{g(b)}^2 \stackrel{\text{(law)}}{=} g(a) B_{f(a)}^2 + \int_a^b dx g'(x) B_{f(x)}^2 \quad (10)$$

where  $f, g : [a, b] \rightarrow \mathbb{R}_+$ , are  $C^1$ , with  $f$  decreasing, and  $g$  increasing. In turn, this implies some extensions of the C-T identities.

Now, it is easily shown that (10) is a particular case of the following Fubini-type identity in law:

$$\int_0^\infty ds \left( \int_0^\infty dB_u \varphi(s, u) \right)^2 \stackrel{\text{(law)}}{=} \int_0^\infty ds \left( \int_0^\infty dB_u \varphi(u, s) \right)^2 \quad (11)$$

where  $\varphi \in L^2([0, \infty[^2; \mathbb{R})$

A striking application of (11) is the following identity in law, obtained with C. Donati-Martin:

$$\int_0^1 ds (B_s - G)^2 \stackrel{\text{(law)}}{=} \int_0^1 ds (B_s - sB_1)^2, \quad (11')$$

where  $G = \int_0^1 du B_u$ .

The general identity (11) is an infinite dimensional extension of the elementary identity in law: if  $\underline{X}_n = (X_1, \dots, X_n)$  is an  $n$ -dimensional sample of the  $N(0, \sigma^2)$  law, then, for any  $n \times n$  matrix:

$$\|A\underline{X}_n\| \stackrel{\text{(law)}}{=} \|A^* \underline{X}_n\|, \quad \text{where } A^* \text{ is the transpose of } A. \quad (12)$$

The above arguments may be developed to give an explanation of the large class of extensions of the C-T identities obtained by Biane [4], between functionals of pairs of diffusions which satisfy a certain duality property.

## 4. Some Limiting Laws of Planar BM

Let  $(Z_t, t \geq 0)$  be a planar BM starting from  $z_0$ , and consider  $(\theta_t, t \geq 0)$  a continuous determination of the argument of  $(Z_u, u \leq t)$  around  $z_1 \neq z_0$ .

Spitzer [24] showed that:

$$\frac{2}{\log t} \theta_t \xrightarrow[t \rightarrow \infty]{\text{(law)}} C_1 \quad (13)$$

where  $C_1$  is a standard Cauchy variable.

This may be refined by decomposing  $\theta_t$  into:  $\theta_t^- + \theta_t^+$ , where:

$$\theta_t^- = \int_0^t d\theta_s 1_{(|Z_s - z_1| \leq R)} \quad \text{and} \quad \theta_t^+ = \int_0^t d\theta_s 1_{(|Z_s - z_1| \geq R)}.$$

Then, considering moreover  $(A_t, t \geq 0)$  an integrable additive functional, we have (Messulam-Yor [16]):

$$\frac{2}{\log t} (\theta_t^-, \theta_t^+, A_t) \xrightarrow[t \rightarrow \infty]{\text{(law)}} (W^-, W^+, c_A \Lambda)$$

where:  $W^- = \int_0^\sigma d\gamma_s 1_{(\beta_s \leq 0)}$ ,  $W^+ = \int_0^\sigma d\gamma_s 1_{(\beta_s \geq 0)}$ ,  $\Lambda = l_\sigma^0$ ,  $c_A$  is a constant depending only on  $A, \beta$  and  $\gamma$  are two independent linear BM's starting from 0,  $\sigma = \inf\{t : \beta_t = 1\}$ , and  $(l_u^0, u \geq 0)$  is the local time of  $\beta$  at 0. The law of  $(W^-, W^+, \Lambda)$  is characterized by:

$$E[\exp(-a\Lambda + ibW^- + icW^+)] = \left( \cosh c + \frac{2a + |b|}{c} \sinh c \right)^{-1} \quad (14)$$

a formula which is very similar to (4), this being easily explained thanks mainly to the Ray-Knight theorem ( $RK_1$ ).

The convergence in law (13) may be further extended by considering  $(\theta_t^1, \dots, \theta_t^n; t \geq 0)$ , the winding numbers of  $(Z_u, u \leq t)$  around  $n$  distinct points  $z_1, \dots, z_n$ .

One obtains (Pitman-Yor [19, 20]):

$$\frac{2}{\log t} (\theta_t^1, \dots, \theta_t^n) \xrightarrow[t \rightarrow \infty]{\text{(law)}} (W_1, \dots, W_n)$$

and the characteristic function of  $(W_1, \dots, W_n)$  is:

$$E\left[\exp i \left( \sum_{j=1}^n \lambda_j W_j \right)\right] = \left( \cosh \left( \sum_{j=1}^n \lambda_j \right) + \frac{\sum |\lambda_j|}{\sum \lambda_j} \sinh \left( \sum_{j=1}^n \lambda_j \right) \right)^{-1}. \quad (15)$$

## 5. Arc sine Laws for Linear Brownian Motion

Let  $(B_t, t \geq 0)$  be the linear BM, starting from 0. Lévy [15] showed that:  $\Gamma_+ \equiv \int_0^\sigma ds 1_{(B_s > 0)}$  follows the arc sine distribution, that is:

$$P(\Gamma_+ \in dt) = \frac{dt}{\pi \sqrt{t(1-t)}}. \quad (16)$$

(Notice that, in paragraph 4, we considered  $\int_0^\sigma ds 1_{(B_s \geq 0)}$ , where  $\sigma \equiv \inf\{t : B_t = 1\}$ ; now,  $\sigma$  is replaced by 1).

The r.v.  $g_1 \equiv \sup\{t < 1 : B_t = 0\}$  is also arc sine distributed, but this is easier to prove.

Here is a recent proof of (16), obtained with M. Barlow and J. Pitman [1], using excursion theory; one can show:

$$\frac{1}{l_t^2} (\Gamma_+(t), \Gamma_-(t)) \stackrel{\text{(law)}}{=} (T_+, T_-), \quad (17)$$

where  $T_+$  and  $T_-$  are two independent stable  $(\frac{1}{2})$  variables.

It now follows that:

$$\Gamma_+ \equiv \Gamma_+(1) \stackrel{\text{(law)}}{=} \frac{T_+}{T_+ + T_-}, \quad \text{which proves (16).}$$

Several infinite dimensional extensions of (17) have now been obtained, jointly with J. Pitman [21]. Here is one: let  $V(t)$  be the infinite sequence of lengths of excursions of  $B$  away from 0, during the time interval  $(0, t)$ , including the last unfinished excursion, arranged in decreasing order, so that:

$$V(t) = (V_1(t), V_2(t), \dots), \quad \text{with } V_1(t) \geq V_2(t) \geq \dots \geq V_n(t) \geq \dots$$

Then, for every  $t > 0$ , and  $s > 0$ ,

$$\frac{V(t)}{t} \stackrel{\text{(law)}}{=} \frac{V(\tau_s)}{\tau_s}, \quad (18)$$

where  $(\tau_s, s \geq 0)$  is the inverse of the local time  $(l_t, t \geq 0)$ .

## 6. Cauchy's Principal Value of Brownian Local Times

Consider again  $(B_t, t \geq 0)$  a linear BM starting from 0,  $(l_t, t \geq 0)$  its local time at 0, and  $(\tau_t, t \geq 0)$  the inverse of  $l$ . As a consequence of the regularity of Brownian local times,

$$H_t \stackrel{\text{def}}{=} \lim_{\varepsilon \rightarrow 0} \int_0^t \frac{ds}{B_s} 1_{(|B_s| \geq \varepsilon)}$$

exists a.s., uniformly in  $t$  in a compact set (whereas:  $\int_0^t \frac{ds}{|B_s|} = \infty$  a.s.). From the scaling property and the Markov property of BM, it follows that:  $(H_{\tau_t}, t \geq 0)$  is a symmetric Cauchy process (with parameter  $\pi!$ ).

Independently, it was remarked by Spitzer [24] that if  $(X_u, u \geq 0)$  is a linear BM independent of  $(\tau_t, t \geq 0)$ , then  $(X_{\tau_t}, t \geq 0)$  is a symmetric Cauchy process.

However, this identity in law does not extend to the two 2-dimensional processes:

$$\left( \frac{1}{\pi} H_{\tau_t}, \tau_t; t \geq 0 \right) \quad \text{and} \quad (X_{\tau_t}, \tau_t; t \geq 0)$$

since we have the formula:

$$E \left[ \exp \left( i \frac{\lambda}{\pi} H_{\tau_t} - \frac{\theta^2}{2} \tau_t \right) \right] = \exp \left( -t \lambda \coth \left( \frac{\lambda}{\theta} \right) \right). \quad (19)$$

From this formula, we deduce:

$$E\left[\exp\left(i\frac{\lambda}{\pi}H_T\right)|l_T=x\right]=\frac{\lambda}{\sinh\lambda}\exp-x(\lambda\coth\lambda-1) \quad (20)$$

where  $T$  is an exponential variable (with parameter  $(\frac{1}{2})$ ) independent of  $B$ . From Lévy's formula (1), we deduce:

$$\left(\frac{1}{\pi}H_T; l_T\right) \stackrel{\text{(law)}}{=} \left(S_1; \frac{1}{2}|Z_1|^2\right) \quad (21)$$

which has not yet received a simple direct explanation.

As a consequence, we obtain, for fixed  $t > 0$ :

$$P(H_t \in dx) = \left(\frac{2}{\pi^3 t}\right)^{1/2} \sum_{n=0}^{\infty} (-1)^n \exp\left(-\left(n + \frac{1}{2}\right)^2 \frac{x^2}{2t}\right) dx. \quad (22)$$

The fact that some relation between the laws of the processes  $(H_t, t \geq 0)$  and  $(S_t, t \geq 0)$  exists may be understood, at least in some sense, via easier identities in law, such as:

$$\int_0^1 \frac{ds}{R_s} \stackrel{\text{(law)}}{=} 2 \left( \int_0^1 ds \tilde{R}_s^2 \right)^{-1/2}, \quad (23)$$

where  $(R_t, t \geq 0)$ , resp:  $(\tilde{R}_t, t \geq 0)$  is a BES process with dimension  $\delta \geq 2$ , resp:  $\tilde{\delta} \equiv 2\delta - 2$ .

All the results presented in this paragraph are taken from Biane-Yor [6]. Principal values of Brownian local times have been studied in depth by Yamada (see, in particular, [27]) and Bertoin [2]; they have also been investigated, for physical purposes, by Ezawa et al [10].

*Note added in proof:* Further results about Cauchy's principal value of local times have now been obtained by Bertoin [30] and Fitzsimmons and Getoor [31].

## References

1. M.T. Barlow, J.W. Pitman, M. Yor: Une extension multidimensionnelle de la loi de l'arc sinus. Sémin. Probab. XXIII. (Lecture Notes in Mathematics, vol. 1372.) Springer, Berlin Heidelberg New York 1989, pp. 294–314
2. J. Bertoin: Applications de la théorie spectrale des cordes vibrantes aux fonctionnelles additives principales d'un brownien réfléchi. Ann. I.H.P., **25**, 3 (1989) 307–323
3. Ph. Biane: Decompositions of Brownian trajectories and some applications. Notes from lectures given at the Probability Winter-school of Wuhan, China, Fall 1990
4. Ph. Biane: Comparaison entre temps d'atteinte et temps de séjour de certaines diffusions réelles. Sémin. Probab. XIX. (Lecture Notes in Mathematics, vol. 1123.) Springer, Berlin Heidelberg New York 1985, pp. 291–296
5. Ph. Biane, M. Yor: Variations sur une formule de Paul Lévy. Ann. I.H.P. **23** (1987) 359–377
6. Ph. Biane, M. Yor: Valeurs principales associées aux temps locaux browniens. Bull. Sci. Maths. **111** (1987) 23–101
7. J.M. Bismut: The Atiyah-Singer theorems. J. Funct. Anal. **57** (1984) 56–99 and 329–348
8. J.M. Bismut: Formules de localisation et formules de Paul Lévy. Astérisque **157–158**. Colloque Paul Lévy sur les processus stochastiques 1988, pp. 37–58

9. Z. Ciesielski, S.J. Taylor: First passage time and sojourn density for Brownian motion in space and the exact Hausdorff measure of the sample path. *Trans. Amer. Math. Soc.* **103** (1962) 434–450
10. H. Ezawa, J.R. Klauder, L.A. Shepp: Vestigial effects of singular potentials in diffusion theory and quantum mechanics. *J. Math. Phys.* **16**, 4 (1975) 783–799
11. T. Jeulin: Ray-Knight's theorems on Brownian local times and Tanaka's formula. Seminar on Stochastic Processes, eds. E. Cinlar, K.L. Chung, R.K. Getoor. Birkhäuser, Basel 1983, pp. 131–142
12. F.B. Knight: Random walks and a sojourn density process of Brownian motion. *Trans. Amer. Math. Soc.* **107** (1963) 56–86
13. J.F. Le Gall, M. Yor: Excursions browniennes et carrés de processus de Bessel. *C. R. Acad. Sci. Paris, Série I* **303** (1986) 73–76
14. P. Lévy: Wiener random functions and other Laplacian random functions. *Proc. Second Berkeley Symp., vol. II*, 1950, pp. 171–186
15. P. Lévy: Sur certains processus stochastiques homogènes. *Comp. Math.* **7** (1939) 283–339.
16. P. Messulam, M. Yor: On D. Williams' “pinching method” and some applications. *J. London Math. Soc.* **26** (1982) 348–364
17. J.W. Pitman: One dimensional Brownian motion and the three-dimensional Bessel process. *Adv. Appl. Prob.* **7** (1975) 511–526
18. J.W. Pitman, M. Yor: A decomposition of Bessel Bridges. *Z. Wahrscheinlichkeitstheorie Verw. Geb.* **59** (1982) 425–457
19. J.W. Pitman, M. Yor: Asymptotic laws of planar Brownian motion. *Ann. Prob.* **14** (1986) 733–779
20. J.W. Pitman, M. Yor: Further asymptotic laws of planar Brownian motion. *Ann. Prob.* **17** (1989) 965–1011
21. J.W. Pitman, M. Yor: Arc sine laws and interval partitions derived from a stable subordinator. Submitted to *Proc. London Math. Soc.* (November 1990)
22. D.B. Ray: Sojourn times of a diffusion process. *III. J. Math.* **7** (1963) 615–630
23. T. Shiga, S. Watanabe: Bessel diffusions as a one-parameter family of diffusion processes. *Z. Wahrscheinlichkeitstheorie Verw. Geb.* **27** (1973) 37–46
24. F. Spitzer: Some theorems concerning 2-dimensional Brownian motion. *Trans. Amer. Math. Soc.* **87** (1958) 187–197
25. D. Williams: On a stopped Brownian motion formula of H.M. Taylor. *Sém. Proba. X. (Lecture Notes in Mathematics, vol. 511.) Springer, Berlin Heidelberg New York 1976*, pp. 235–239
26. D. Williams: Path decomposition and continuity of local time for one-dimensional diffusions I. *Proc. London Math. Soc.* **28** (3) (1974) 738–768
27. T. Yamada: On some limit theorems for occupation times of one-dimensional Brownian motion and its continuous additive functionals locally of zero energy. *J. Math. Kyoto Univ.* **26** (2) (1986) 309–322
28. M. Yor: Remarques sur une formule de Paul Lévy. *Sém. Proba. XIV. (Lecture Notes in Mathematics, vol. 784.) Springer, Berlin Heidelberg New York 1980*, pp. 343–346
29. M. Yor: Une explication du théorème de Ciesielski-Taylor. To appear in *Ann. I.H.P.* **27** (2) (1991)
30. J. Bertoin: Complements on the Hilbert transform and the fractional derivatives of Brownian local times. *J. Math. Kyoto* **30** (4) (1990) 651–670
31. P.J. Fitzsimmons and R.K. Getoor: On the distribution of the Hilbert transform of the local time of a symmetric Lévy process. To appear in *Ann. Prob.*

# The Stability of Minkowski Spacetime

*Demetrios Christodoulou*

Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street  
New York, NY 10012, USA

According to Einstein, spacetime is a 4-dimensional manifold  $M$  with a metric  $g_{\mu\nu}$  of signature (3, 1). General relativity is a unified theory of space, time and gravitation in which the connection of  $g_{\mu\nu}$  is identified with the gravitational force and the Einstein equations

$$R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R = 8\pi T_{\mu\nu}$$

hold, where  $R_{\mu\nu}$  and  $R$  are respectively the Ricci curvature and scalar curvature of  $g_{\mu\nu}$  and  $T_{\mu\nu}$  is the energy tensor of matter. In the absence of matter the equations reduce to the Einstein vacuum equations for the spacetime manifold:

$$R_{\mu\nu} = 0.$$

These equations are at first sight a degenerate differential system. That is, the null space of the symbol  $\sigma_\xi$  at a given covector  $\xi$  is nonzero for all covectors  $\xi$ . This is due to the fact that the equations are generally covariant; proper account must be taken of the geometric equivalence of metrics related by diffeomorphism. This is done by considering  $\sigma_\xi$  not on the space of 2-covariant symmetric tensors at a point but rather on the quotient of this space by the equivalence relation induced by the symbol of the diffeomorphisms; two such tensors are equivalent if they differ by  $\xi X + X \xi$  for some covector  $X$ . The null space of  $\sigma_\xi$  is then found to be nonzero if and only if  $\xi$  belongs to the null cone defined by the metric  $g$ . The Einstein equations are therefore of hyperbolic character.

The central mathematical problem of the theory is the initial value problem. An initial data set is a 3-dimensional manifold  $\Sigma$  with a positive definite metric  $\bar{g}_{ij}$  and a 2-covariant symmetric tensorfield  $k_{ij}$ . The problem is to find a 4-dimensional manifold  $M$  with a metric  $g_{\mu\nu}$  of signature (3, 1) satisfying the Einstein vacuum equations and an imbedding of  $\Sigma$  into  $M$  such that  $\bar{g}_{ij}$  and  $k_{ij}$  are respectively the first and second fundamental forms induced on  $\Sigma$  by the imbedding. The Einstein vacuum equations impose on the initial data set the constraint equations:

$$\begin{aligned}\bar{\nabla}^i k_{ij} - \bar{\nabla}_j \operatorname{tr} k &= 0 \\ \bar{R} - |k|^2 + (\operatorname{tr} k)^2 &= 0\end{aligned}$$

where  $\bar{R}$  is the scalar curvature of  $\bar{g}_{ij}$ . These are respectively the Codazzi and Gauss equations of the imbedding of  $\Sigma$  in  $M$ .

The local existence of solutions of the initial value problem was proven by Choquet-Bruhat [1] by introducing a harmonic system of coordinates, thereby reducing the equations to a differential system for the metric components in these coordinates which is hyperbolic in the standard sense. By a global solution of the initial value problem we mean a spacetime  $(M, g_{\mu\nu})$  which is geodesically complete. Penrose [2] found a basic obstruction to global existence in his singularity theorem: the spacetime cannot be geodesically complete if  $\Sigma$  is non-compact and contains a trapped sphere, namely a surface  $S$  diffeomorphic to  $S^2$  such that the divergence of the outgoing future directed null normals to  $S$  in  $M$  is everywhere negative.

An initial data set is called asymptotically flat if the complement of a compact set in  $\Sigma$  is diffeomorphic to the complement of a ball in  $\mathbf{R}^3$ ,  $\bar{g}_{ij}$  is a complete metric on  $\Sigma$  and the curvature of  $\bar{g}_{ij}$  as well as  $k_{ij}$  tend to zero at infinity in an appropriate way. Since the formulation of general relativity, the Minkowski spacetime has been the only known global solution of the vacuum Einstein equations arising from asymptotically flat initial data. A basic problem in the theory is the question of stability of Minkowski spacetime, that is, whether any asymptotically flat initial data set which is sufficiently close to the trivial one gives rise to a global solution of the vacuum Einstein equations. In a recent joint work Sergiu Klainerman and myself have answered this question in the affirmative. We have also studied in detail the asymptotic behaviour of the solutions.

Before giving an outline of our proof I wish to give an indication of the difficulty involved. The naive approach to the problem would be to try to extend the solution obtained by Choquet-Bruhat of the system of equations for the metric components in harmonic coordinates, to a global one. It turns out that this approach would work if the space dimension were greater than 3. However, it fails in 3 space dimensions, a fact already recognized by Choquet-Bruhat.

Our proof uses two main ideas. The first is the relationship between conserved quantities and symmetry and the second is the relationship between symmetry and the causal structure of spacetime. The first idea goes back to Noether. Consider a field theory in a given spacetime whose field equations are derivable from an action  $S$ . The energy tensor is then defined by:

$$T_{\mu\nu} = \frac{\delta S}{\delta g^{\mu\nu}}$$

and the invariance of  $S$  under diffeomorphisms implies that  $T_{\mu\nu}$  is divergence-free:

$$\nabla^\mu T_{\mu\nu} = 0.$$

Now suppose that  $X$  is a vectorfield generating a 1-parameter group of isometries of  $(M, g)$  (killing vectorfield). Then the 1-form

$$P_\mu = -T_{\mu\nu}X^\nu$$

is divergence-free:

$$\nabla^\mu P_\mu = 0.$$

It follows that the integral

$$\int_{\Sigma} *P$$

on a Cauchy hypersurface  $\Sigma$  is a conserved quantity, that is, its value is the same for all Cauchy hypersurfaces. If the action is invariant under conformal transformations of the metric then the energy tensor is trace-free  $\text{tr } T = 0$  and the above considerations extend to the case where  $X$  generates a 1-parameter group of conformal isometries of  $(M, g)$  (conformal killing vectorfield). An important requirement on a physical theory is that the energy tensor should satisfy the positivity condition

$$T(X_1, X_2) \geq 0$$

for any two future directed timelike vectors  $X_1, X_2$  at a point. If the vectorfield  $X$  above is timelike and future directed then the quantity

$$\int_{\Sigma} *P = \int_{\Sigma} T(X, N) d\mu_{\bar{g}}$$

is nonnegative,  $N$  being the future directed unit normal to  $\Sigma$ . As its value is the same as that on the Cauchy hypersurface on which the initial data is given, it provides an estimate for the solution in terms of the initial data.

In the case of gravitation the energy tensor, defined as above,

$$G_{\mu\nu} = \frac{\delta S}{\delta g^{\mu\nu}}$$

vanishes, as this expresses the Euler-Lagrange equations of gravitation, namely the Einstein equations. Thus the above considerations fail at first sight. The way out of this impasse is to consider the Bianchi identities

$$\nabla_{[\alpha} R_{\beta\gamma]\delta\varepsilon} = 0$$

(where  $[ ]$  stands for cyclic permutation). We define a Weyl field  $W_{\alpha\beta\gamma\delta}$ , in a given spacetime, to be a 4-covariant tensorfield possessing the algebraic symmetries of the Weyl or conformal curvature tensor. The natural field equations for a Weyl field are the Bianchi equations

$$\nabla_{[\alpha} W_{\beta\gamma]\delta\varepsilon} = 0$$

which we write simply as

$$DW = 0.$$

In a spacetime satisfying the vacuum Einstein equations the curvature is an example of a Weyl field satisfying the Bianchi equations. In a 4-dimensional spacetime the dual  $*W$  of a Weyl field  $W$  is also a Weyl field and if  $W$  satisfies the Bianchi equations so does  $*W$ . The operator  $D$  although formally identical to the exterior derivative, is not an exterior differential operator and  $D^2 \neq 0$ . As a consequence, the Bianchi equations imply an algebraic condition:

$$R_{\mu}^{\alpha\beta\gamma} * W_{\nu\alpha\beta\gamma} - R_{\nu}^{\alpha\beta\gamma} * W_{\mu\alpha\beta\gamma} = 0.$$

The Bianchi equations are conformally covariant. If  $f$  is a conformal isometry of

$(M, g)$ , that is  $f * g = \Omega^2 g$  for some positive function  $\Omega$ , and  $W$  is a solution of the Bianchi equations then so is  $\Omega^{-1} f * W$ .

Associated to a Weyl field  $W$  is a 4-covariant tensorfield  $Q$ , quadratic in  $W$ , called the Bel-Robinson tensor [4]:

$$Q_{\alpha\beta\gamma\delta} = W_{\alpha\beta\gamma\sigma} W_{\beta}{}^{\sigma}{}_{\delta} + *W_{\alpha\beta\gamma\sigma} * W_{\beta}{}^{\sigma}{}_{\delta}.$$

It is totally symmetric and trace-free and satisfies the following positivity condition

$$Q(X_1, X_2, X_3, X_4) \geq 0$$

for any four future directed timelike vectors  $X_1, X_2, X_3, X_4$  at a point, with equality if and only if  $W$  vanishes at that point. Furthermore, if  $W$  satisfies the Bianchi equations then  $Q$  is divergence free:

$$\nabla^\alpha Q_{\alpha\beta\gamma\delta} = 0.$$

It follows that if  $X_1, X_2, X_3$  are three vectorfields each generating a 1-parameter group of conformal isometries of  $(M, g)$ , then the 1-form

$$P = -Q(\cdot, X_1, X_2, X_3)$$

is divergence-free, consequently the integral

$$\int_{\Sigma} *P$$

on a Cauchy hypersurface  $\Sigma$  is a conserved quantity, which is positive definite in the case that  $X_1, X_2, X_3$  are all timelike and future directed.

Given a Weyl field  $W$  and a vectorfield  $X$  the Lie derivative  $\mathcal{L}_X W$  of  $W$  with respect to  $X$  is not in general a Weyl field. However we can define a modified Lie derivative  $\hat{\mathcal{L}}_X W$  which is a Weyl field:

$$\begin{aligned} \hat{\mathcal{L}}_X W_{\alpha\beta\gamma\delta} &= \mathcal{L}_X W_{\alpha\beta\gamma\delta} - \frac{1}{8} \operatorname{tr} \pi W_{\alpha\beta\gamma\delta} \\ &\quad - \frac{1}{2} (\hat{\pi}_{\mu\alpha} W_{\mu\beta\gamma\delta} + \hat{\pi}_{\mu\beta} W_{\alpha\mu\gamma\delta} \\ &\quad + \hat{\pi}_{\mu\gamma} W_{\alpha\beta\mu\delta} + \hat{\pi}_{\mu\delta} W_{\alpha\beta\gamma\mu}) \end{aligned}$$

where  $\pi_{\mu\nu} = \mathcal{L}_X g_{\mu\nu}$  and  $\hat{\pi}$  is the deformation tensor of  $X$ , namely the trace-free part of  $\pi$ . The modified Lie derivative commutes with the Hodge dual:

$$\hat{\mathcal{L}}_X * W = * \hat{\mathcal{L}}_X W.$$

As a consequence of the linearity and the conformal covariance of the Bianchi equations, if  $W$  is a solution of these equations and  $X$  is a vectorfield generating a 1-parameter group of conformal isometries  $f_t$ , then

$$\hat{\mathcal{L}}_X W = \frac{d}{dt} (\Omega_t^{-1} f_t * W)|_{t=0}$$

is also a solution of the same equations. Therefore the above considerations regarding conserved quantities can be applied to the Weyl field  $\hat{\mathcal{L}}_X W$  as well.

Minkowski spacetime has a 15 dimensional conformal group  $O(4, 2)$  consisting of the abelian subgroups of translations and inverted translations, the scalings and the Lorentz group  $O(3, 1)$ . The generators of time translations and inverted time translations are the only timelike conformal killing vectorfields. A general spacetime will not have a nontrivial conformal group. If the condition that the vectorfields in the above considerations are conformal killing is dropped, then, although the quantities

$$\int_{\Sigma} *P$$

will not be conserved, their growth shall be determined by the spacetime integrals of expressions which are quadratic in  $W$  and linear in  $\pi$ , the deformation tensors of the vectorfields. The idea is to find a subgroup of  $O(4, 2)$  and an action of this subgroup on the spacetime manifold having the following property: the deformation tensors of the vectorfields generating the action should decay at infinity in such a way that the growth of the corresponding quantities is bounded in terms of the quantities themselves. It turns out that the subgroup consisting of time translations, scalings, inverted time translations and the rotation group  $O(3)$ , leaving the total energy vector invariant, suffices.

The problem is then how to define the action of these groups on a general spacetime arising from asymptotically flat initial data in such a way as to satisfy the above requirement. Now the group of time translations is the simplest to define and is related to the choice of a time function  $t$ . As our argument is one of continuity starting from the initial Cauchy hypersurface, a time function seems to enter the problem naturally. A canonical choice is that of a maximal time function, namely one whose level sets  $\Sigma_t$  are maximal spacelike hypersurfaces. The second fundamental form  $k_{ij}$  of  $\Sigma_t$  is then trace-free:  $\text{tr } k = 0$ . There is a unique such function  $t$  with the property that the total momentum vector relative to  $\Sigma_t$ , that is, the projection to  $\Sigma_t$  of the total energy vector, vanishes. To define the action of the other groups we consider the fact that the spacetime is expected to be asymptotically flat. Since these groups act canonically on Minkowski spacetime, there is a canonical action defined at infinity. What we need is a way to extend this action to the spacetime. This is provided by the causal structure. The causal structure on a manifold  $M$  is the fundamental structure defined on  $M$  by a metric of signature  $(\dim M - 1, 1)$ . The causal future  $J^+(S)$  of a set  $S \subset M$  is the set of points  $q$  which can be reached by a future directed causal curve initiating at  $S$ . Similarly  $J^-(S)$  is the set of points  $q$  which can be reached, from  $S$ , by a past directed causal curve. The specification of  $J^+(p)$  and  $J^-(p)$  for every  $p \in M$  defines the causal structure, which is equivalent to the conformal geometry of  $M$ . In our problem we consider the boundaries of the causal pasts of a 1-parameter family of surfaces at “infinity” which are related to each other by a time translation. However we cannot quite start from “infinity” since the existence of a global solution is precisely what we wish to establish. Nevertheless, following our continuity argument we can assume that a spacetime slab has been constructed. The role of “infinity” is then played by the final maximal hypersurface  $\Sigma_{t_*}$ . We construct a 1-parameter family of surfaces diffeomorphic to  $S^2$  on  $\Sigma_{t_*}$  by solving a certain equation of motion for 2-surfaces

on a 3-dimensional manifold. We then consider the inner boundaries of the causal pasts of these surfaces in the spacetime slab. These null hypersurfaces  $C_u$  are the level sets of what we call the optical function  $u$ . Let  $S_{t,u}$  be the surfaces of intersection of the  $C_u$  with the  $\Sigma_t$ . Let  $l$  and  $\underline{l}$  be respectively the outgoing and incoming null normals to  $S_{t,u}$  whose component along  $T$ , the generator of time translations, is equal to  $T$ . Then we have

$$T = \frac{1}{2}(l + \underline{l})$$

and we define the generator of scalings by

$$S = \frac{1}{2}(u\underline{l} + u\underline{l})$$

and the generator of inverted time translations by

$$K = \frac{1}{2}(u^2 l + u^2 \underline{l})$$

where

$$u = u + 2r, \quad r = \sqrt{\frac{A}{4\pi}}$$

and  $A$  is the area of  $S_{t,u}$ . An action of the rotation group  $O(3)$  on  $\Sigma_{t,*}$  is defined starting from the standard action on the sphere at infinity in such a way that the group orbits are the level surfaces of  $u$  on  $\Sigma_{t,*}$ . The action is then extended to the spacetime slab by conjugation: Given an element  $O \in O(3)$  and a point  $p \in S_{t,u}$ , to obtain the point  $Op$  we follow the generator of  $C_u$  through  $p$  toward the future until  $p_*$ , the point of intersection with  $\Sigma_{t,*}$ . We then move to  $Op_*$  and follow the generator of  $C_u$  through that point toward the past until the point of intersection with  $\Sigma_t$ . This point, which again lies on  $S_{t,u}$ , is the sought for point  $Op$ . The three rotation vectorfields  ${}^{(a)}\Omega$ ,  $a = 1, 2, 3$ , generating this action satisfy:

$$[{}^{(a)}\Omega, l] = 0,$$

$$g({}^{(a)}\Omega, l) = g({}^{(a)}\Omega, T) = 0$$

and, of course, the commutation relations of the Lie algebra of  $O(3)$ :

$$[{}^{(a)}\Omega, {}^{(b)}\Omega] = \epsilon_{abc} {}^{(c)}\Omega.$$

The group orbits are the surfaces  $S_{t,u}$ .

By the above construction, the deformation tensors of the generating vectorfields depend entirely on the geometric properties of the hypersurfaces  $C_u$  and  $\Sigma_t$ . These properties differ significantly from those in the case of Minkowski spacetime. Consider for example, on a given  $C_u$ , the second fundamental form  $\theta$  of  $S_{t,u}$  relative to  $\Sigma_t$ , in particular the ratio

$$f = \frac{|\hat{\theta}|^2}{(\text{tr } \theta)^2}$$

where  $\hat{\theta}$  is the trace-free part of  $\theta$ . This ratio, in contrast to the case in Minkowski spacetime, does not tend to zero as  $t \rightarrow \infty$ . In fact,  $\lim_{t \rightarrow \infty} f$  measures the flux of energy radiated to infinity. Another example is the area of  $S_{t,u}$  on a given  $C_u$ , which

in Minkowski spacetime verifies  $r - t = O(1)$  as  $t \rightarrow \infty$ , while in general we find

$$r - t = -2M_0 \log t + O(1),$$

where  $M_0$  is the initial total mass. These differences, which are absent in the case of space dimension greater than 3, are the reason for the failure of the naive approach I mentioned earlier.

In carrying out the plan outlined above we encounter a certain difficulty. The positive mass theorem of Schoen-Yau [5] and Witten [6] implies that the curvature of any nontrivial initial data set decays not faster than  $d^{-3}$  where  $d$  is the distance from a given point. However, general initial data with  $d^{-3}$  decay lead to a logarithmic divergence in the construction of the optical function. The difficulty is overcome by observing that the leading terms at infinity in the initial data correspond to the Schwarzschild solution which is both spherically symmetric and static and is therefore annihilated by both the rotations and the time translations. We therefore take as our basic Weyl fields not the spacetime curvature  $R$  but rather  $\hat{\mathcal{L}}_T R$  and  $\hat{\mathcal{L}}_0 R$ , where 0 stands for the collection  $\{^{(a)}\Omega : a = 1, 2, 3\}$ . Due to the fact that there are no nontrivial solutions of the vacuum Einstein equations which have a Cauchy hypersurface diffeomorphic to  $\mathbf{R}^3$  and are either spherically symmetric or stationary, we are able to control  $R$  itself in terms of  $\hat{\mathcal{L}}_0 R$  and  $\hat{\mathcal{L}}_T R$ .

The proof of the theorem is by the method of continuity and it involves a bootstrap argument. Using an appropriate version of the local existence theorem we can assume that the spacetime is maximally extended up to a value  $t_*$  of the time function. This value is defined to be the maximal one such that certain geometric quantities defined by the level sets of the time function and the optical function remain bounded by a small positive number  $\varepsilon_0$ . These quantities control in particular the isoperimetric constant of the surfaces  $S_{t,u}$ , on which the Sobolev inequalities depend. It then follows that a certain norm of the deformation tensors of the vectorfields  $T, S, K$  and 0 in the spacetime slab bounded by  $\Sigma_0$  and  $\Sigma_{t_*}$  is less than another small positive number  $\varepsilon_1$ . We then consider the 1-form  $P$ , where

$$\begin{aligned} P &= P_1 + P_2, \\ P_1 &= -Q(\hat{\mathcal{L}}_0 R)(\cdot, \bar{K}, \bar{K}, T) - Q(\hat{\mathcal{L}}_T R)(\cdot, \bar{K}, \bar{K}, \bar{K}), \\ P_2 &= -Q(\hat{\mathcal{L}}_0^2 R)(\cdot, \bar{K}, \bar{K}, T) - Q(\hat{\mathcal{L}}_0 \hat{\mathcal{L}}_T R)(\cdot, \bar{K}, \bar{K}, \bar{K}) \\ &\quad - Q(\hat{\mathcal{L}}_S \hat{\mathcal{L}}_T R)(\cdot, \bar{K}, \bar{K}, \bar{K}) - Q(\hat{\mathcal{L}}_T^2 R)(\cdot, \bar{K}, \bar{K}, \bar{K}) \end{aligned}$$

and

$$\bar{K} = K + T.$$

We define the quantity  $E = \max\{E_1, E_2\}$  where

$$E_1 = \sup_t \int_{\Sigma_t} *P, \quad E_2 = \sup_u \int_{C_u} *P$$

and everything is restricted to the spacetime slab under consideration. The crucial point is the estimate of the error terms which control the growth of  $\int_{\Sigma_t} *P$  and  $\int_{C_u} *P$ . Using the bound for the deformation tensors we are able to estimate the integral in

the spacetime slab of these error terms by  $c\varepsilon_1 E$  and thus arrive at an inequality of the form

$$E \leq c(D + \varepsilon_1 E)$$

where  $D$  stands for initial data. When  $\varepsilon_1$  is chosen sufficiently small, which is achieved by choosing  $\varepsilon_0$  suitably small, this implies  $E \leq cD$ . On the other hand we are able to show that the aforementioned geometric quantities associated to the level sets of  $t$  and  $u$  are bounded by  $cE$ . Thus if  $D$  is suitably small this bound does not exceed  $\varepsilon_0/2$ , which by continuity contradicts the maximality of  $t_*$ , unless of course  $t_* = \infty$ , in which case the theorem is proved.

The estimate of the error terms would fail if it were not for the fact that the worst error terms vanish due to a simple algebraic identity: if  $A, B, C$  are any three symmetric trace-free 2-dimensional matrices then  $\text{tr}(ABC) = 0$ . The reason why such matrices appear here can be traced back to the symbol of the Einstein equations. As I mentioned at the beginning the null space of  $\sigma_\xi$  is nontrivial if and only if  $\xi$  is a null convector. But we can say more; when  $\xi$  is a (nonzero) null covector, the null space of  $\sigma_\xi$  is in fact isomorphic to the space of symmetric trace-free 2-dimensional matrices. This is therefore the space of dynamical degrees of freedom of the gravitational field at a point. There is no product in this space because for any two such matrices  $A, B$  we have  $AB + BA - I \text{tr}(AB) = 0$ . This implies the identity mentioned above.

It remains for me to state the smallness condition on the initial data which is required. Take a point  $p \in \Sigma_0 = \Sigma$  and a positive real number  $a$ . Let  $d_p$  be the distance function on  $\Sigma$  from  $p$ . Set:

$$\begin{aligned} D(p, a) = \sup_{\Sigma} & \{a^{-2}(d_p^2 + a^2)^3 |\bar{Ric}|^2\} \\ & + a^{-3} \left\{ \int_{\Sigma} \sum_{l=0}^3 (d_p^2 + a^2)^{l+1} |\bar{\nabla}^l k|^2 d\mu_{\bar{g}} \right. \\ & \left. + \int_{\Sigma} \sum_{l=0} (d_p^2 + a^2)^{l+3} |\bar{\nabla}^l B|^2 d\mu_{\bar{g}} \right\} \end{aligned}$$

here  $\bar{Ric}$  is the Ricci curvature and  $B$  the Bach tensor or conformal curvature of  $(\Sigma, \bar{g})$ :

$$B = \text{curl } \hat{\bar{Ric}},$$

where  $\hat{\bar{Ric}}$  is the trace-free part of  $\bar{Ric}$ . Then it is the dimensionless invariant

$$\inf_{p \in \Sigma, a > 0} D(p, a)$$

which must be sufficiently small.

In concluding I would like to emphasize that many deeper and more difficult mathematical problems remain in general relativity. Among these are the formation of trapped surfaces, the nature of the singularities and the so called cosmic censorship question: whether or not singularities are generically proceeded by a region of trapped surfaces. All these are aspects of the initial value problem in the large.

## References

1. Y. Choquet-Bruhat: Théorème d'existence pour certains systèmes d'équations aux dérivées partielles non linéaires. *Acta Math.* **88** (1952) 141–225
2. R. Penrose: Gravitational collapse and spacetime singularities. *Phys. Rev. Lett.* **14** (1965) 57–59
3. D. Christodoulou, S. Klainerman: The global nonlinear stability of the Minkowski space. *Ann. Math Studies* (to appear)
4. L. Bel: Introduction d'un tenseur du quatrième ordre. *C. R. Acad. Sci. Paris* **248** (1959) 1094–1096
5. R. Schoen, S.T. Yau: On the proof of the positive mass conjecture in general relativity. *Commun. Math. Phys.* **65** (1979) 45–76
6. E. Witten: A new proof of the positive energy theorem. *Commun. Math. Phys.* **80** (1981) 381–402



# Harmonic Maps with Values into Spheres

Jean-Michel Coron

Université Paris-Sud, Département de Mathématiques  
Bâtiment 425, F-91405 Orsay, France

## 0. Introduction

Let  $(M, g)$  and  $(N, h)$  be two compact Riemannian manifolds of dimension  $m$  and  $n$ ; the manifold  $M$  may have a boundary but not  $N$ . Without any loss of generality we may assume, using Nash's theorem, that  $N$  is isometrically embedded in  $\mathbb{R}^k$ . Associated with a map  $u$  from  $M$  into  $N$  is the Dirichlet energy

$$E(u) = \int_M e(u)(x) dM(x) \quad (0.1)$$

where  $e(u)(x)$  is the square of the Hilbert-Schmidt norm of  $u'(x) : T_x M \rightarrow T_{u(x)} N$ . Let us note that  $E$  is defined on the Sobolev space  $H^1(M; N) = \{u \in H^1(M; \mathbb{R}^k); u(x) \in N \text{ for a.e. } x\}$ . The critical point of  $E$  on  $H^1(M; N)$  are called harmonic maps. The Euler-Lagrange equation satisfied by a harmonic map is

$$\Delta_M u^i = g^{\alpha\beta} A_{u(x)}^i \left( \frac{\partial u}{\partial x^\alpha}, \frac{\partial u}{\partial x^\beta} \right) \quad \text{for } i \in \{1, \dots, m\} \quad (0.2)$$

where  $A_u$  is the second fundamental form of  $N$ . This generalization of the usual harmonic functions is due to J. Eells and J. Sampson [ES]. Since their pioneer work, the subject of harmonic maps between manifolds has drawn the attention of many analysts, geometers and physicists. For quite complete surveys on that subject we refer to two papers by J. Eells and L. Lemaire: [EL1] and [EL2].

When  $N$  is the unit sphere  $S^n$  of  $\mathbb{R}^{n+1}$  and  $M = \bar{\Omega}$ , where  $\Omega$  is a smooth bounded open set of  $\mathbb{R}^m$ , the energy of  $u$  is (with standard notations)

$$E(u) = \int_\Omega |\nabla u|^2 = \int_\Omega u_i^\alpha u_i^\alpha \quad (0.3)$$

and  $u : M \rightarrow N$  is harmonic if and only if  $u \in H^1(M; N)$  and satisfies, in the sense of distributions,

$$-\Delta u = u |\nabla u|^2. \quad (0.4)$$

A first natural question is the regularity of harmonic maps in the interior of  $M$ . One easily checks that if  $m = 1$  any harmonic map is smooth. This is an open question for  $m = 2$  but quite interesting partial results are known (Section 1). For  $m \geq 3$  a harmonic map needs not to be smooth on the interior of  $M$ .

(e.g.  $u : B^m = \{x \in \mathbb{R}^m; |x| < 1\} \rightarrow S^{m-1}$ ,  $u(x) = x/|x|$ ). On the other hand S. Hildebrandt, H. Kaul and K.O. Widman have proved in [HKW] that if the image of  $u$  is included in a “small” ball then  $u$  is smooth; for  $N = S^n$  a ball included in an open hemisphere is “small”; in that case their result is optimal since  $u : B^3 \rightarrow S^3$  defined by  $u(x) = (x/|x|, 0)$  is harmonic. Unfortunately there is no general result on the singular set  $S(u)$  of a harmonic map  $u$ . Much more is known on  $S(u)$  if  $u$  is a *minimizing* harmonic map. A map  $u : M \rightarrow N$  is a minimizing harmonic map if for any  $v : M \rightarrow N$  such that  $v = u$  on  $\partial M$  then  $E(u) \leq E(v)$ . In that case R. Schoen and K. Uhlenbeck have proved in [SU1] and [SU2] (see also M. Giaquinta and E. Giusti [GG] if the image of  $u$  is included in a chart) that  $S(u)$  is a compact set of dimension less or equal to  $(m - 3)$  and that, near a singular point,  $u$  behaves like a singular homogeneous minimizing harmonic map from  $B^m$  into  $N$ . Such maps are called minimizing tangent maps (MTM). R. Schoen and K. Uhlenbeck have proved in [SU3] that there are no MTM from  $B^m$  into  $S^n$  if  $m \leq d(n)$  where  $d(3) := 3$  and  $d(n) := 1 + \min\{n/2, 5\}$  otherwise. A classification of the MTM for  $m = 3$  and  $N = S^2$  is given in [BCL] – see Section 2 –. Moreover in [BCL] a sharp lower bound of  $E(u)$  for  $u : \Omega \subset \mathbb{R}^3 \rightarrow S^2$  with prescribed singularities is given (Section 3). In Section 4 we give examples of MTM due to F.H. Lin [L] and [CG] for  $N = S^n$ .

When  $n = m - 1$  and  $N = S^m$  then if the degree of  $u$  restricted to  $\partial M$  is not zero,  $S(u)$  cannot be empty. On the other hand R. Hardt and F.H. Lin have proved in [HL1] that  $S(u)$  may be not empty even if  $u$  restricted to  $\partial\Omega$  is of degree zero. More precisely they have in particular constructed smooth maps  $\gamma : \partial B^3 \rightarrow S^2$  of degree zero such that (gap phenomenon)

$$\text{Min}\{E(v); v \in H_\gamma^1(B^3; S^2)\} < \text{Inf}\{E(v); v \in H_\gamma^1(B^3; S^2) \cap C^1(\bar{B}^3; S^2)\} \quad (0.5)$$

where  $H_\gamma^1(B^3; S^2) = \{u \in H^1(B^3; S^2); u = \gamma \text{ on } \partial B^3\}$ . Inequality (0.5) implies that  $H_\gamma^1(B^3; S^2) \cap C^\infty(B^3; S^2)$  is not dense in  $H_\gamma^1(B^3; S^2)$ . More generally one may ask under which condition  $C^\infty(M; N)$  is dense in the Sobolev space  $W^{1,p}(M; N)$ . In Section 5 we will describe recent results concerning this problem.

Inequality (0.5) leads also to the question: is the infimum in the right hand side of (0.5) achieved? The answer is not known but it is proved in [BBC] that a positive answer to that question is equivalent to the smoothness of a minimizer for a relaxed energy associated to  $E$ . Moreover, using various relaxed energies, it is possible to prove – see [BBC] – that, if the degree of  $\gamma$  is not zero or if (0.5) holds, then the Dirichlet problem with boundary data  $\gamma$  has infinitely many solutions.

Our paper is organised as follows:

1. Regularity of Harmonic Maps from a Surface.
2. Classification of MTM for  $m = 3$  and  $N = S^2$ .
3. Lower Bound of  $E(u)$  for  $\Omega \subset \mathbb{R}^3$ ,  $N = S^2$  when  $S(u)$  is Known.
4. Example of MTM for  $N = S^n$ .
5. Gap Phenomenon and the Problem of the Density of Smooth Maps.
6. Relaxed Energies.

## 1. Regularity of Harmonic Maps from a Surface

In this section we assume  $m = 2$ . Let us first mention that the following question is still open

**Question 1.1.** *Let  $u$  be a harmonic map from  $M$  into  $N$ . Is  $u$  smooth in the interior of  $M$ ?*

A positive answer has been given under some extra assumptions. One has

**Theorem 1.2.** *Assume  $u : M \rightarrow N$  is harmonic; then  $u$  is smooth in the following cases:*

- (i)  $u$  is a minimizing harmonic map (C.B. Morrey [M]);
- (ii) there exists  $x_0$  in the interior of  $M$  such that  $u$  is smooth on  $M \setminus \{x_0\}$  (J. Sacks and K. Uhlenbeck [SaU]);
- (iii)  $\phi := (|u_1|^2 - |u_2|^2 - 2i < u_1, u_2 >) (dx_1 + idx_2)^2$  is holomorphic (M. Grüter [G] if  $\phi = 0$ , R. Schoen [Sc] in the general case);
- (iv)  $N$  is a sphere (F. Hélein [H2]).

Case (iv) is a quite recent result. Let us briefly sketch the proof of (iv). F. Hélein first notes that

$$\nabla u^i = u^j (\nabla u^i - u^i \nabla u^j), \quad \text{for } 1 \leq i \leq 3. \quad (1.1)$$

Next, using (0.4), one has (see [C, KRS and Sh]):

$$\operatorname{div}(u^j \nabla u^i - u^i \nabla u^j) = 0 \quad \text{for } 1 \leq i \leq 3 \text{ and } 1 \leq j \leq 3; \quad (1.2)$$

hence there exists  $B^{ij}$  such that (at least locally)

$$\operatorname{curl} B^{ij} = u^j \nabla u^i - u^i \nabla u^j \quad \text{for } 1 \leq i \leq 3 \text{ and } 1 \leq j \leq 3. \quad (1.3)$$

From (1.2) and (1.3) F. Hélein gets

$$\Delta u^i = \det(\nabla u^j, \nabla B^{ij}) \quad \text{for } 1 \leq i \leq 3. \quad (1.4)$$

Finally the continuity of  $u^i$  – hence the smoothness of  $u$  (see [LU] Chap. 8) – follows from (1.4) and a result due to H. Wente [We].

We end up this section by mentioning two recent results slightly related to Question 1.1 or Theorem 1.2.

**Theorem 1.3** (F. Hélein [H1]). *Let  $u$  be a quasi-conformal homeomorphism between two Riemannian surfaces. Then, if  $(|u_1|^2 - |u_2|^2 - 2i < u_1, u_2 >) (dx_1 + idx_2)^2$  is holomorphic,  $u$  is a smooth harmonic map.*

**Theorem 1.4** [CH]. *Let  $u$  be a smooth harmonic map from  $M$  into  $N$  and let  $\varphi : M \rightarrow M$  be a continuous map homotopic to the identity map, then  $E(u) \leq E(u \circ \varphi)$ .*

## 2. Classification of the MTM for $m = 3$ and $N = S^2$

In this section  $m = 3$  and  $N = S^2$ . Then a MTM  $u$  is a non-constant minimizing harmonic map from  $\bar{B}^3$  into  $S^2$  such that  $u(x) = \omega(x/|x|)$  for some harmonic map  $\omega : S^2 \rightarrow S^2$ . Let us recall that these maps are important since they give the behavior of a minimizing harmonic map from a 3-dimensional manifold into  $S^2$  (see [SU1, Si and GW]).

The classification of MTM is given by the following theorem proved in [BCL].

**Theorem 2.1.** *The map  $u : \bar{B}^3 \rightarrow S^2$  is a MTM if there exists  $R$  in  $SO(3)$  such that  $u(x) = \pm Rx/|x|$ .*

*Sketch of the proof of Theorem 2.1.* Let  $u$  be a MTM from  $\bar{B}^3$  into  $S^2$ . We have  $u(x) = \omega(x/|x|)$  where  $\omega$  is a harmonic map from  $S^2$  into  $S^2$  and therefore there exist two polynomials  $P$  and  $Q$  prime together such that, if we denote by  $\pi$  the stereographic projection  $\pi : S^2 \rightarrow (\mathbb{R}^2 \cup \{\infty\} \times \{0\}) \simeq \mathbb{C} \cup \{\infty\}$  with pole  $(0,0,1)$

$$\text{either } \pi\omega\pi^{-1}(z) = P(z)/Q(z) \quad \text{or} \quad \pi\omega\pi^{-1}(z) = P(\bar{z})/Q(\bar{z}) \quad (2.1)$$

– see [EL1 (10.6)] –. One first notices that, since  $u$  is a minimizing harmonic map, then for any smooth vector field  $X$  on  $B^3$  with compact support  $\frac{d}{de}E(u(x + eX(x)))|_{e=0} = 0$  which leads to

$$\mathbb{R}^3 \ni \int_{S^2} S|\nabla\omega|^2 dS = 0. \quad (2.2)$$

Then one checks that, if  $d := \text{Max}(\deg P, \deg Q) = 1$ , (2.2) implies that  $\omega = \pm R$  for some  $R$  in  $SO(3)$ . In the case  $|d| \geq 2$  it is proved in [BCL] that one can decrease the energy of  $u$  by splitting the singularity  $\{0\}$  into  $|d|$  distinct points.

It remains only to verify that  $Rx/|x|$  is a MTM. There are now a lot of proofs available for that – see [BCL, L and CG] –. The shortest method is perhaps to say that since there are singular harmonic maps  $\bar{B}^3 \rightarrow S^2$  there exists at least a MTM from  $\bar{B}^3$  into  $S^2$ ; such a MTM needs to be of the form  $\pm R_0 x/|x|$  for some  $R_0$  in  $SO(3)$ ; but clearly, if  $\pm R_0 x/|x|$  is a MTM, then  $\pm Rx/|x|$  is also a MTM for any  $R$  in  $SO(3)$ .

**Remark 2.2.** In an earlier paper [HKL] R. Hardt, D. Kinderlehrer and F.H. Lin had proved that the degree  $d$  of a MTM is bounded by a universal constant. Moreover, R. Cohen, R. Hardt, D. Kinderlehrer, S.Y. Lin and M. Luskin had proved numerically that if  $\pi\omega\pi^{-1}(z) = z^2$ ,  $u$  is not a MTM ([HKLL]; see also some recent numerical studies by F. Alouges [A]).

### 3. Lower Bound of $E(u)$ for $\Omega \subset \mathbb{R}^3$ and $N = S^2$ when $S(u)$ is Known

We start with the special case  $\Omega = \mathbb{R}^3$ . Let  $(a_i)_{1 \leq i \leq \ell}$  be  $\ell$  points in  $\mathbb{R}^3$  and  $(d_i)_{1 \leq i \leq \ell}$  be  $\ell$  integers. We are interested in the value of

$$I = \inf \left\{ \int_{\mathbb{R}^3} |\nabla u|^2 ; u \in C^1(\mathbb{R}^3 \setminus \{a_1, \dots, a_\ell\}; S^2) \right. \\ \left. \text{and } \deg(u, a_i) = d_i \quad \text{for } 1 \leq i \leq \ell \right\}. \quad (3.1)$$

In (3.1)  $\deg(u, a_i)$  denotes the degree of  $u$  restricted to a small ball centered at  $a_i$ . One easily checks that, if  $m < \infty$ , then

$$\sum_{i=1}^{\ell} d_i = 0. \quad (3.2)$$

So we will assume that (3.2) holds true. In order to give an explicit formula for  $I$ , we construct with the points  $a_i$  such that  $d_i > 0$  a family of points  $(P_j)_{1 \leq j \leq p}$  (called “positive” points) in the following way : each  $P_j$  belongs to the set  $\{a_i ; d_i > 0\}$  and each  $a_i$  with  $d_i > 0$  is repeated exactly  $d_i$  times in the family  $(P_j)$ ; we do the same with the set  $\{a_i ; d_i < 0\}$  and get the family of “negative” points  $(N_j)_{1 \leq j \leq n}$ . It follows from (3.2) that  $p = n$ . Let  $L$  be defined by

$$L = \text{Min} \left\{ \sum_{j=1}^p |P_j - N_{\sigma(j)}| ; \sigma \text{ is a permutation of } \{1, 2, \dots, p\} \right\}. \quad (3.3)$$

Then we have [BCL]

**Theorem 3.1.** *The infimum in (3.1) is never achieved and  $I = 8\pi L$ .*

*Sketch of the proof of  $I = 8\pi L$ .* a)  $I \leq 8\pi L$ . This part relies on the following lemma (dipole construction).

**Lemma 3.2.** *Let  $P$  and  $N$  be two points in  $\mathbb{R}^3$  and let  $\varepsilon > 0$ . There exists a sequence of maps  $(u_n)_n$  in  $C^1(\mathbb{R}^3 \setminus \{P, N\}; S^2)$  such that*

$$E(u_n) \rightarrow 8\pi|P - N| \text{ as } n \rightarrow \infty \\ \deg(u_n, P) = 1, \deg(u_n, N) = -1 \\ u_n = \text{North pole on } \{Q \in \mathbb{R}^3 ; \text{dist}(Q, [P, N]) \geq \varepsilon\}.$$

Lemma 3.2 gives  $I \leq 8\pi L$  if  $n = p = 1$ . The general case is obtained by gluing together dipoles for the pairs  $(P_i, N_{\sigma(i)})$ .

b)  $I \geq 8\pi L$ . Two different proofs are known.

α) See [BCL]. Let, for  $u \in \{v \in C^1(\mathbb{R}^3 \setminus \{a_1, \dots, a_\ell\}; S^2) ; \deg(v, a_i) = d_i\}$ ,

$$D(u) = (u \cdot (u_2 \times u_3), u \cdot (u_3 \times u_1), u \cdot (u_1 \times u_2)).$$

Then one has

$$\operatorname{div} D(u) = 4\pi \left( \sum_{i=1}^n (\delta_{P_i} - \delta_{N_i}) \right) \quad \text{and} \quad |\nabla u|^2 \geq 2|D(u)|. \quad (3.4)$$

From (3.4) one gets

$$\begin{aligned} E(u) \geq 8\pi \operatorname{Min} \left\{ \sum_{i=1}^n (\xi(P_i) - \xi(N_i)); \xi : \mathbb{R}^3 \rightarrow \mathbb{R} \quad \text{such that} \right. \\ \left. |\xi(x) - \xi(y)| \leq |x - y| \quad \forall x \in \mathbb{R}^3, \forall y \in \mathbb{R}^3 \right\}. \end{aligned} \quad (3.5)$$

Finally, using Kantorovitch's theorem on min-max and the characterization due to G. Birkhoff of the extremal points of the set of doubly stochastic matrices, one can prove that the right-hand side in (3.5) is  $8\pi L$ .

$\beta$ ) This proof is due to F. Almgren, W. Browder and E.H. Lieb [ABL]. It relies on the co-area formula and goes as follows. Let  $u$  be in  $C^1(\mathbb{R}^3 \setminus \{a_1, \dots, a_\ell\}; S^2)$  with  $\deg(u, a_i) = d_i$  for any  $i$  in  $\{1, \dots, \ell\}$ . One has

$$E(u) \geq 2 \int_{\mathbb{R}^3} J_2(u) dx = 2 \int_{S^2} \mathcal{H}^1(u^{-1}(S)) d\sigma(S), \quad (3.6)$$

$$\partial(u^{-1}(S)) = \sum_{i=1}^P (\delta_{P_i} - \delta_{N_i}). \quad (3.7)$$

Finally from (3.6) and (3.7) one gets  $E(u) \geq 8\pi L$ .

Two extensions of Theorem 3.1 to the case  $\Omega \neq \mathbb{R}^3$  are possible. We assume that the points  $a_i$  are in  $\Omega$ . Let

$$I_1 = \operatorname{Inf}\{E(u); u \in C^1(\Omega \setminus \{a_1, \dots, a_\ell\}; S^2),$$

$$\deg(u, a_i) = d_i \text{ for } 1 \leq i \leq \ell \text{ and } u \text{ is constant on } \partial\Omega\}$$

$$I_2 = \operatorname{Inf}\{E(u); u \in C^1(\Omega \setminus \{a_1, \dots, a_\ell\}; S^2), \deg(u, a_i) = d_i \text{ for } 1 \leq i \leq \ell\}.$$

Let us denote by  $D$  the geodesic distance in  $\Omega$ . Let

$$L_2 = \operatorname{Min} \left\{ \sum_{i=1}^p D(P_i, N_{\sigma(i)}); \sigma \text{ is a permutation of } \{1, \dots, p\} \right\}$$

if  $p = n$ ;  $L_2 = +\infty$  if  $p \neq n$ . Let us define  $L_1$  for  $p \leq n$  by

$$\begin{aligned} L_1 = \operatorname{Min} \left\{ \sum_{i=1}^p D(P_i, N_{\sigma(i)}) + \sum_{i=p+1}^n \operatorname{dist}(N_{\sigma(i)}, \partial\Omega); \right. \\ \left. \sigma \text{ is a permutation of } \{1, \dots, n\} \right\}. \end{aligned} \quad (3.8)$$

For  $p \geq n$   $L_1$  is defined by replacing in (3.8)  $n$  by  $p$ ,  $p$  by  $n$ ,  $N$  by  $P$  and  $P$  by  $N$ . Then one has (see [BCL])  $I_1 = 8\pi L_1$  and  $I_2 = 8\pi L_2$ .

**Remark 3.3.** a) The analogous of Theorems 3.1 for the liquid crystals energy has been obtained by M. Giaquinta, G. Modica and J. Souček in [GMS3]. b) F. Almgren, E.H. Lieb in [AL] and R. Hardt, F.H. Lin in [HL2] have obtained estimates on the number of singularities of minimizing harmonic maps from  $\Omega \subset B^3$  into  $S^2$ .

## 4. Examples of MTM

The following maps have been proved to be MTM:

- a)  $u_0 : B^{n+m} \subset \mathbb{R}^{n+m} = \mathbb{R}^{n+1} \times \mathbb{R}^{m-1} \rightarrow S^n$ ,  $u_0(x', x'') = x'/|x'|$ .
- b)  $u_0 : B^{2n} \rightarrow S^n$ ,  $u_0(x) = H(x/|x|)$  where  $H : S^{2n-1} \rightarrow S^n$  are the Hopf maps related to the multiplication of complex numbers ( $n = 2$ ), quaternions ( $n = 4$ ) and Cayley numbers ( $n = 8$ ). Example a) for  $m = 1$  is due to F.H. Lin [L]. Example a) for  $m \geq 2$  and Example b) are proved in [CG].

Lin's proof relies on the null Lagrangian method. It can be divided into three steps.

*Step 1.* Let  $u : \mathbb{R}^{n+1} \rightarrow S^n$  then

$$|\nabla u|^2 \geq (n-1)^{-1} \{\text{tr}(\nabla u)^2 - (\text{div } u)^2\} \text{ with equality if } u = u_0. \quad (4.1)$$

*Step 2.* Let  $u : B^{n+1} \rightarrow \mathbb{R}^{n+1}$  then

$$\text{tr}(\nabla u)^2 - (\text{div } u)^2 = \text{div}\{(\text{div } u)u - u \cdot \nabla u\}. \quad (4.2)$$

*Step 3.* Let  $u \in H^1(B^{n+1}; S^n)$  with  $u = u_0$  on  $\partial B^{n+1}$ . Inequality (4.1) gives

$$\int_{B^{n+1}} |\nabla u|^2 \geq (n-1)^{-1} \int_{B^{n+1}} (\text{tr}(\nabla u)^2 - (\text{div } u)^2). \quad (4.3)$$

Using (4.2) one gets

$$\int_{B^{n+1}} (\text{tr } u)^2 - (\text{div } u)^2 = \int_{B^{n+1}} (\text{tr } u_0)^2 - (\text{div } u_0)^2. \quad (4.4)$$

Finally using (4.3), (4.4) and (4.1) for  $u = u_0$  one has  $\int_{B^{n+1}} |\nabla u|^2 \geq \int_{B^{n+1}} |\nabla u_0|^2$ .

The proof in [CG] relies on the co-area formula as in [ABL] and on a projection-averaging procedure. Let us explain the method for Example b) with the quaternionic Hopf map. Let  $u : B^8 \rightarrow S^4$  such that  $u = u_0$  on  $\partial B^8$ . We want to prove that

$$E(u) \geq E(u_0). \quad (4.5)$$

For  $\psi$  in  $G_3(\mathbb{R}^5)$ , write  $\pi_\psi : S^4 \rightarrow S^2$  for the nearest point projection. One first checks that

$$E(u) = \frac{2}{3} \oint_{G_3(\mathbb{R}^5)} E(\pi_\psi \circ u). \quad (4.6)$$

Let  $v = \pi_\psi \circ u$  and  $v_0 = \pi_\psi \circ u_0$ . Inequality (4.5) is a consequence of (4.6) and

$$E(v) \geq E(v_0). \quad (4.7)$$

The proof of (4.7) goes as follow

$$E(v) \geq 2 \int_{B^8} J_2(v) = 2 \int_{S^2} \mathcal{H}^6(v^{-1}(S)) = \int_{S^2} \mathcal{H}^6(v^{-1}(-S) \cup v^{-1}(S)) ; \quad (4.8)$$

Note that, at least if  $v$  is not too singular – see [CG] for more details –,

$$\begin{aligned} \partial(v^{-1}(-S) \cup v^{-1}(S)) &= (\pi_\psi \circ H)^{-1}(S) - (\pi_\psi \circ H)^{-1}(-S) \\ &= \partial(v_0^{-1}(-S) \cup v_0^{-1}(S)). \end{aligned} \quad (4.9)$$

One can prove that, for some complex structure on  $\mathbb{R}^8$ ,  $v_0^{-1}(-S) \cup v_0^{-1}(S)$  is a complex variety; hence, by (4.9) and [F, p. 435 and 652],

$$\mathcal{H}^6(v^{-1}(-S) \cup v^{-1}(S)) \geq \mathcal{H}^6(v_0^{-1}(-S) \cup v_0^{-1}(S)). \quad (4.10)$$

Inequality (4.7) follows from (4.8), (4.10) and

$$E(v_0) = 2 \int_{B^8} J_2(v_0) = \int_{S^2} \mathcal{H}^6(v_0^{-1}(-S) \cup v_0^{-1}(S)). \quad (4.11)$$

## 5. The Gap Phenomenon and the Problem of the Density of Smooth Maps

Let us start with a quite interesting theorem proved by R. Hardt and F.H. Lin in [HL1].

**Theorem 5.1** (Gap Phenomenon). *The exist smooth maps  $\gamma : \partial B^3 \rightarrow S^2$  of degree zero such that*

$$\begin{aligned} \text{Min}\{E(u); u \in H^1(B^3; S^2)\} \quad \text{and} \quad u = \gamma \text{ on } \partial B^3 \\ < \text{Inf}\{E(u); u \in C^1(\bar{B}^3; S^2)\} \quad \text{and} \quad u = \gamma \text{ on } \partial B^3 \}. \end{aligned} \quad (5.1)$$

*Sketch of the construction of  $\gamma$ .* We follow a method proposed by H. Brezis in [Br]. Let  $\varepsilon$  be a small positive number. Let  $P_1 = (0, 0, 1 - \varepsilon)$ ,  $N_1 = (0, 0, 1 + \varepsilon)$ ,  $N_2 = (0, 0, -1 + \varepsilon)$ ,  $P_2 = (0, 0, -1 - \varepsilon)$ . By Lemma 3.2 there exists a map  $u_\varepsilon : \mathbb{R}^3 \rightarrow S^2$  smooth on  $\mathbb{R}^3 \setminus \{N_1, N_2, P_1, P_2\}$  such that the degree of  $u_\varepsilon$  at the  $N_i$  (resp.  $P_i$ ) is  $-1$  (resp.  $+1$ ) and

$$\int_{\mathbb{R}^3} |\nabla u_\varepsilon|^2 \leq 8\pi(2\varepsilon + 2\varepsilon) + \varepsilon. \quad (5.2)$$

One takes for  $\gamma$  the restriction of  $u_\varepsilon$  to  $\partial B^3$ . The degree of  $\gamma$  is zero. For  $v$  in  $C^1(\bar{B}^3; S^2)$  with  $v = \gamma$  on  $\partial B^3$  let  $w : \mathbb{R}^3 \rightarrow S^2$  be defined by  $w = v$  on  $B^3$  and  $w = u_\varepsilon$  on  $\mathbb{R}^3 \setminus B^3$ . That map is continuous on  $\mathbb{R}^3 \setminus \{N_1, P_2\}$  and has degree  $+1$  at  $P_2$ ,  $-1$  at  $N_1$  hence by Theorem 3.1

$$\int_{\mathbb{R}^3} |\nabla w|^2 \geq 8\pi(2 + 2\varepsilon) \quad (5.3)$$

which implies, with (5.2),

$$\int_{B^3} |\nabla v|^2 \geq 16\pi(1 - \varepsilon) - \varepsilon. \quad (5.4)$$

The gap phenomenon follows from (5.2) and (5.4).

Note that the gap phenomenon implies that even if  $\gamma$  is of degree zero  $C_\gamma^1(B^3; S^2) = \{v; v \in C^1(B^3; S^2) \text{ and } v = \gamma \text{ on } \partial B^3\}$  may not be dense in  $H_\gamma^1(B^3; S^2) = \{v \in H^1(B^3; S^2); v = \gamma \text{ on } \partial B^3\}$ . Similar non-density had been observed previously by R. Schoen and K. Uhlenbeck; they had proved in [SU2]

**Theorem 5.2.** *The map  $u_0$  in  $H^1(B^3; S^2)$  defined by  $u_0(x) = x/|x|$  cannot be approximated in the  $H^1$ -norm by maps in  $C^1(\bar{B}^3; S^2)$ .*

*Proof of Theorem 5.2.* Assume  $u^n \in C^1(\bar{B}^3; S^2)$  satisfy  $u^n \rightarrow u_0$  in  $H^1(B^3; \mathbb{R}^3)$ . This implies that  $D(u^n) \rightarrow D(u_0)$  in  $L^1(B^3; \mathbb{R}^3)$  but  $\operatorname{div} D(u^n) = 0$  and  $\operatorname{div} D(u_0) = 4\pi\delta_0$ . A contradiction.

This leads to the question raised by J. Eells and L. Lemaire in [EL2]: let  $p$  be in  $[1, \infty)$ , is  $C^\infty(M, N)$  dense in the Sobolev space  $W^{1,p}(M; N)$ ? The answer is now known. One has

**Theorem 5.3.** *If  $p \geq m$ , then  $C^\infty(M; N)$  is dense in  $W^{1,p}(M; N)$ . If  $p < m$  then  $C^\infty(M; N)$  is dense in  $W^{1,p}(M; N)$  if and only if the homotopy group  $\pi_{[p]}(N)$  is trivial ( $[p]$  is the largest integer less or equal to  $p$ ).*

Theorem 5.3 is easy if  $p > m$  since in that case, by the Sobolev embeddings,  $W^{1,p}(M; N) \subset C(M; N)$ . The case  $p = m$  has been proved by R. Schoen and K. Uhlenbeck in [SU2] and [SU3]. When  $p < m$  and  $\pi_{[p]}(N) \neq 0$ , F. Bethuel and X. Zheng have constructed in [BZ] a map in  $W^{1,p}(M; N)$  which cannot be approximated in the  $W^{1,p}$ -norm by maps in  $C^\infty(M; N)$ ; their proof uses previous arguments due to B. White [Wh]. The fact that  $\pi_{[p]}(N)$  is trivial implies density is a difficult theorem due to F. Bethuel [B2].

When one does not have density we may ask for a characterization of the closure  $\widehat{W^{1,p}}(M, N)$  of  $C^\infty(M, N)$  in  $W^{1,p}(M; N)$ . Only partial results are known. One has

**Theorem 5.4.** *Assume that  $N$  is  $([p] - 1)$  connected,  $H_{[p]}(N)$  is torsion free and that  $\pi_1(N)$  is abelian if  $[p] = 1$ . Then the two following conditions are equivalent*

$$u \in \widehat{W^{1,p}}(M, N) \quad (5.5)$$

*the pullback by  $u$  of any closed  $[p]$ -form on  $N$  is (weakly) closed.* (5.6)

The implication (5.5)  $\Rightarrow$  (5.6) has been noticed by R. Schoen and K. Uhlenbeck in [SU3] – it holds without any topological assumption on  $N$  –. The converse has been proved by F. Bethuel in [B1] for  $p = 2$ ,  $M = B^3$  and  $N = S^2$ , by F. Demengel in [D] for  $1 \leq p < 2$ ,  $M = B^m$  and  $N = S^1$ ; the general case is proved in [BCDH].

Theorem 5.3 leads also to the question: which size of singularities one has to allow in order to have density ? F. Bethuel has proved in [B2] – see also [BZ] for  $M = B^m$ ,  $N = S^{m-1}$ ,  $2 \leq p < m$ ,

**Theorem 5.5.** *Assume  $p < m$ . The maps from  $M$  into  $N$  which are smooth except on a manifold of dimension  $m - [p] - 1$  are dense in  $W^{1,p}(M; N)$ .*

## 6. Relaxed Energies

The gap phenomenon leads naturally to the following question

**Question 6.1.** *Let  $\gamma : \partial B^3 \rightarrow S^2$  be a smooth map of degree zero. Is the infimum  $\inf\{E(u); u \in C^1(\bar{B}^3; S^2) \text{ and } u = \gamma \text{ on } \partial B^3\}$  achieved?*

The answer to this question is still not known but in [BBC] it has been proved that a positive answer to that question is equivalent to the regularity of a minimizer for a relaxed energy  $E_1$  associated to  $E$ . In order to define  $E_1$  we need some notations. Let  $H_\gamma^1(B^3; S^2) = \{u \in H^1(B^3; S^2); u = \gamma \text{ on } \partial B^3\}$ ,  $C_\gamma^1(B^3; S^2) = \{u \in C^1(\bar{B}^3; S^2); u = \gamma \text{ on } \partial B^3\}$ ,  $R_\gamma(B^3; S^2) = \{u \in H_\gamma^1(B^3; S^2); u \text{ is } C^1 \text{ except at a finite number of points}\}$ . Recall – see Theorem 5.4 – ([B2] or [BZ]) that  $R_\gamma$  is dense in  $H_\gamma^1(B^3; S^2)$  for the  $H^1$ -norm. For  $u$  in  $R_\gamma$  let

$$L(u) = \min \left\{ \sum_{i=1}^n |P_i - N_{\sigma(i)}| ; \quad \sigma \text{ is a permutation of } \{1, \dots, n\} \right\} \quad (6.1)$$

where the  $P_i$  (resp.  $N_i$ ) are the singularities of  $u$  of positive (resp. negative) degree counted according to their degree. By a result of [BCL], for any  $u$  in  $R_\gamma$ ,

$$L(u) = \frac{1}{4\pi} \sup \left\{ \int_{B^3} D(u) \cdot \nabla \theta - \int_{\partial B^3} \theta \operatorname{Jac}(\gamma); |\theta(x) - \theta(y)| \leq |x - y| \forall x, \forall y \right\}$$

from which it follows that  $L(u)$  makes sense for  $u$  in  $H_\gamma^1(B^3; S^2)$ . Let now  $E_1 : H_\gamma^1(B^3; S^2) \rightarrow [0, +\infty]$  be defined by

$$E_1(u) = E(u) + 8\pi L(u).$$

Note  $E_1 = E$  on  $C_\gamma^1(B^3; S^2)$ . The following properties of  $E_1$  are proved in [BBC].

$$E_1 \text{ is weakly lower semicontinuous} \quad (6.2)$$

$$\min \{E_1(u); u \in H_\gamma^1(B^3; S^2)\} = \inf \{E(u); u \in C_\gamma^1(B^3; S^2)\} \quad (6.3)$$

$$E_1(u) = \inf \{\liminf E(u^n); u^n \in C_\gamma^1(B^3; S^2) \text{ and } u^n \rightarrow u \text{ a.e.}\}. \quad (6.4)$$

In particular a positive answer to question: is any minimizer of  $E_1$  smooth ? will give a positive answer to Question 6.1. Recently M. Giaquinta, G. Modica and J. Souček have proved in [GMS2] that the (eventual) singular set of any minimizer of  $E_1$  has Hausdorff dimension less or equal to 1.

Even if this approach does not give, for the moment being, the answer to Question 6.1 it allows to prove

**Theorem 6.2.** Let  $\gamma : \partial B^3 \rightarrow S^2$  be a smooth map. Assume that either the degree of  $\gamma$  is not zero or the degree of  $\gamma$  is zero and (5.1) holds. Then there are infinitely many harmonic maps in  $H_\gamma^1(B^3; S^2)$ .

Let us give the main steps of the proof of Theorem 6.2 when the degree of  $\gamma$  is zero and (5.1) holds. Let for  $\lambda$  in  $[0, 1]$ ,  $E_\lambda(u) = E(u) + 8\pi\lambda L(u)$ . From (6.2) we get

$$E_\lambda \text{ is weakly lower semicontinuous;} \quad (6.5)$$

hence

$$\text{the infimum } m_\lambda \text{ of } E_\lambda \text{ on } H_\gamma^1(B^3; S^2) \text{ is achieved.} \quad (6.6)$$

Using the gap phenomenon one can prove, see [BBC],

$$m_\lambda > m_0 \quad \text{for any } \lambda > 0. \quad (6.7)$$

Note also that clearly

$$\lim_{\lambda \rightarrow 0} m_\lambda = m_0. \quad (6.8)$$

The infinitely many harmonic maps in  $H_\gamma^1(B^3; S^2)$  follows from (6.6), (6.7) and (6.8) and

$$\text{any minimizer of } E_\lambda \text{ is a harmonic map.} \quad (6.9)$$

It remains to prove (6.9). Let  $u$  be in  $R_\gamma$  and let  $\varphi$  be in  $C_0^1(B^3; \mathbb{R}^3)$ ; then the singularities of  $(u + \varepsilon\varphi)/|u + \varepsilon\varphi|$  are the same as those of  $(u + \varepsilon\varphi)/|u + \varepsilon\varphi|$  if  $|\varepsilon||\varphi|_{L^\infty} < 1$  hence, see (6.1),

$$L((u + \varepsilon\varphi)/|u + \varepsilon\varphi|) = L(u) \quad \text{if } |\varepsilon||\varphi|_{L^\infty} < 1. \quad (6.10)$$

By the density of  $R_\gamma$  in  $H_\gamma^1(B^3; S^2)$ , (6.10) also holds for  $u$  in  $H_\gamma^1(B^3; S^2)$ ; it gives (6.9).

**Remark 6.3.** a) F. Bethuel and H. Brezis have proved in [BB] that any minimizer of  $E_\lambda$  for  $\lambda \in [0, 1]$  is smooth on  $\bar{B}^3$  except at a finite number of points. b) A quite interesting different approach to relaxed energies have been given by M. Giaquinta, G. Modica and J. Souček in [GMS1] and [GMS2].

**Note added in proof.** In a very interesting paper (to appear in C.R. Acad. Sci. Paris) F. Hélein has given a positive answer to Question 1.1.

## References

- [A] Alouges, F.: Un schéma numérique pour le calcul d'applications harmoniques de  $\mathbb{R}^3$  dans la sphère. C.R. Acad. Sci. Paris, 1990
- [ABL] Almgren, F., Browder, W., Lieb, E.H.: Co-area, liquid crystals and minimal surfaces. (Lecture Notes in Mathematics, vol. 1306.) Springer, Berlin Heidelberg New York 1988, pp. 1–12
- [AL] Almgren, F., Lieb, E.H.: Singularities of energy minimizing maps from the ball to the sphere: examples, counter-examples and bounds. Ann. Math. 128 (1988) 483–530

- [B1] Bethuel, F.: A characterization of maps in  $H^1(B^3, S^2)$  which can be approximated by smooth maps. Ann. IHP Analyse non Linéaire (to appear)
- [B2] Bethuel, F.: Approximations dans des espaces de Sobolev et groupes d'homotopie. C.R. Acad. Sci. Paris **307** (1988) 293–296, and The approximation problem for Sobolev maps between two manifolds. Acta Math. (to appear)
- [BB] Bethuel, F., Brezis, H.: Régularité de minima pour des énergies relaxées. C.R. Acad. Sci. Paris (1990), and Regularity of minimizers of relaxed problems for harmonic maps. Preprint
- [BBC] Bethuel, F., Brezis, H., Coron, J.M.: Relaxed energies for harmonic maps. In: Variational problems, H. Berestycki, J.M. Coron, I. Ekeland (eds.). Birkhäuser, 1990
- [BCDH] Bethuel, F., Coron, J.M., Demengel, F., Hélein, F.: A cohomological criterion for density of smooth maps in Sobolev spaces between two manifolds. In: Defects, singularities and patterns in nematic liquid crystals, J.M. Coron, J.M. Ghidaglia, F. Hélein (eds.), Kluwer Academic Press, 1991
- [Br] Brezis, H.:  $S^k$ -valued maps with singularities. In: Topics in calculus of variations, M. Giaquinta (ed.). (Lecture Notes in Mathematics, vol. 1365.) Springer, Berlin Heidelberg New York 1989, pp. 1–30
- [BCL] Brezis, H., Coron, J.M., Lieb, E.H.: Harmonic maps with defects. Comm. Math. Phys. **107** (1986) 649–705
- [BZ] Bethuel, F., Zheng, X.: Density of smooth functions between two manifolds in Sobolev spaces. J. Funct. Anal. **80** (1988) 60–75
- [C] Chen, Y.: Weak solutions to the evolutions problems of harmonic maps. Math. Z. **201** (1989) 69–74
- [CG] Coron, J.M., Gulliver, R.: Minimizing  $p$ -harmonic maps into spheres. J. Reine Angew. Math. **401** (1989) 82–100
- [CH] Coron, J.M., Hélein, F.: Minimizing harmonic diffeomorphisms. Compositio Math. **69** (1989) 175–228
- [D] Demengel, F.: Une caractérisation des applications de  $W^{1,p}(B^N, S^1)$  qui ne peuvent être approchées par des fonctions  $C^\infty$ . C.R. Acad. Sci. Paris **310** (1990) 553–557
- [EL1] Eells, J., Lemaire, L.: Report on harmonic maps. Bull. London Math. Soc. **10** (1978) 1–68
- [EL2] Eells, J., Lemaire, L.: Another report on harmonic maps. Bull. London Math. Soc. **20** (1988) 385–524
- [ES] Eells, J., Sampson, J.H.: Harmonic mappings of Riemannian manifolds. Amer. J. Math. **86** (1964) 109–160
- [F] Federer, H.: Geometric measure theory. Springer, Berlin Heidelberg New-York, 1969
- [G] Grüter, M.: Regularity of weak  $H$ -surfaces. J. Reine Angew. Math. **329** (1981) 1–15
- [GG] Giaquinta, M., Giusti, E.: The singular set of the minima of certain quadratic functionals. Ann. S.N.S. Pisa **11** (1984) 45–55
- [GMS1] Giaquinta, M., Modica, G., Souček, G.: Cartesian currents, weak diffeomorphisms and existence theorems in nonlinear elasticity. Arch. Rational Mech. Anal. **106** (1989) 97–159
- [GMS2] Giaquinta, M., Modica, G., Souček, G.: The Dirichlet energy of mappings with values into the sphere. Manuscripta Math. **65** (1989) 489–507
- [GMS3] Giaquinta, M., Modica, G., Souček, G.: Liquid crystals: relaxed energies, dipoles, singular lines and singular points. Annali S.N.S. Pisa (to appear)
- [GW] Gulliver, R., White, B.: The rate of a convergence of a harmonic map at a singular point. Math. Ann. **283** (1989) 539–549

- [H1] Hélein, F.: Homéomorphismes quasi conformes entre surfaces riemanniennes. *C.R. Acad. Sci. Paris* **307** (1988) 725–730
- [H2] Hélein, F.: Régularité des applications faiblement harmoniques entre une surface et une sphère. *C.R. Acad. Sci. Paris* **311** (1990) 519–524
- [HKL] Hardt, R., Kinderlehrer, D., Lin, F.H.: Existence and partial regularity of static liquid crystal configurations. *Comm. Math. Phys.* **105** (1986) 547–570
- [HKLL] Hardt, R., Kinderlehrer, D., Lin, S.Y., Luskin, M.: Minimum energy configurations for liquid crystals: computational results. In: *Theory and applications of liquid crystals*, J.L. Ericksen, D. Kinderlehrer (eds.), IMA **5**, 1987
- [HKW] Hildebrandt, S., Kaul, H., Widman, K.O.: An existence theorem for harmonic mappings of Riemannian manifolds. *Acta Math.* **138** (1977) 1–16
- [HL1] Hardt, R., Lin, F.H.: A remark on  $H^1$  mappings. *Manuscripta Math.* **56** (1986) 1–10
- [HL2] Hardt, R., Lin, F.H.: Stabilities of singularities of minimizing harmonic maps. *J. Diff. Geom.* **29** (1989) 113–123
- [KRS] Keller, J., Rubinstein, J., Sternberg, P.: Reaction-diffusion processes and evolution to harmonic maps. Preprint
- [L] Lin, F.H.: Une remarque sur l'application  $x/|x|$ . *C.R. Acad. Sci. Paris* **305** (1987) 529–531
- [LU] Ladyshenskaya, O.A., Vral'tseva: Linear and quasilinear elliptic equations. Academic Press, New York 1968
- [M] Morrey Jr., C.B.: The problem of plateau on a Riemannian manifold. *Ann. Math.* **49** (1948) 807–851
- [SaU] Sacks, J., Uhlenbeck, K.: The existence of minimal immersions of 2-spheres. *Ann. Math.* **113** (1981) 1–24
- [Sc] Schoen, R.: Analytic aspects of the harmonic maps problem. *M.S.R.I.* **2** (1984) 321–358
- [Sh] Shatah, J.: Weak solutions and development of singularities of the  $SU(2)$   $\sigma$ -model. *Comm. Pure Appl. Math.* **41** (1988) 459–469
- [Si] Simon, L.: Asymptotics for a class of nonlinear evolution equations with applications to geometric problems. *Ann. Math.* **118** (1983) 525–571
- [SU1] Schoen, R., Uhlenbeck, K.: A regularity theory for harmonic maps. *J. Diff. Geom.* **17** (1982) 307–335
- [SU2] Schoen, R., Uhlenbeck, K.: Boundary regularity and the Dirichlet problem for harmonic maps. *J. Diff. Geom.* (1983) 253–268
- [SU3] Schoen, R., Uhlenbeck, K.: Regularity of minimizing harmonic maps into the sphere. *Invent. math.* **78** (1984) 89–100
- [SU4] Schoen, R., Uhlenbeck, K.: Approximation theorems for Sobolev mappings. Preprint
- [We] Wente, H.: An existence theorem for surfaces of constant mean curvature. *J. Math. Anal. Appl.* **26** (1969) 318–344
- [Wh] White, B.: Infima of energy functionals in homotopy classes. *J. Diff. Geom.* **23** (1986) 127–142



# Isometric Embeddings of Riemannian Manifolds

Matthias Günther

Universität Leipzig, Sektion Mathematik, Augustusplatz 10  
D-7010 Leipzig, Fed. Rep. of Germany

## 1. Explanation of the Problem

Let  $M$  be an  $n$ -dimensional manifold of class  $C^\infty$  and  $g$  any given Riemannian metric on  $M$ . We will consider the following classical problem motivated by differential geometry. Does there exist an embedding  $u = (u^1, \dots, u^q) : M \rightarrow \mathbb{R}^q$  such that the usual euclidian metric of  $\mathbb{R}^q$  induces on the submanifold  $u(M)$  the given metric  $g$ ? In other words,  $u$  must satisfy

$$F(u) := du \cdot du = g, \quad (1)$$

or in local coordinates

$$\sum_{l=1}^q \frac{\partial u^l}{\partial x^i} \frac{\partial u^l}{\partial x^j} = g_{ij}.$$

The dot in (1) denotes the usual scalar product of  $\mathbb{R}^q$ . The notion embedding means, that  $u$  is locally an immersion and globally a homeomorphism of  $M$  onto the subspace  $u(M)$  of  $\mathbb{R}^q$ . If an embedding  $u : M \rightarrow \mathbb{R}^q$  satisfies (1) on the whole  $M$ , we speak of an *isometric embedding*. If  $u$  is an immersion and a solution of (1) in a (possibly small) neighbourhood of any point of  $M$ , we speak of a *local isometric embedding*. A further question is the regularity of the embedding in dependence on the regularity of the metric. And finally, what can be said about the minimal value of  $q$ ?

We will give some historical remarks. There exists a great number of beautiful papers which handle the isometric embedding problem (local or global) under further assumptions on the manifold  $M$  or the metric  $g$  (e.g. special values of dimension  $n$ , positivity assumptions on the curvature), but we are interested only in the general problem. Janet (1926), Cartan (1927) and Burstin (1931) proved the existence of a local isometric embedding with  $q = n(n + 1)/2$  in the analytical case; the essentials of their proofs are suitable applications of the Cauchy-Kowalewski theorem. Up to date a corresponding result (with the same value of  $q$  and without further assumptions) is unknown in the nonanalytical case, even for the dimension  $n = 2$ . Nash (1954) and Kuiper (1955) proved the existence of a (global) isometric embedding of class  $C^1$  with  $n < q \leq 2n + 1$  (in fact their results are more subtle),

provided that the metric  $g$  is continuous. In an outstanding paper Nash (1956) showed the existence of a (global) isometric embedding  $u \in C^s(M, \mathbb{R}^q)$  ( $s \geq 3, s = \infty$ ) if  $g \in C^s$  and  $q = 3n(n+1)/2 + 4n$  in the compact case,  $q = (n+1)q_{\text{com}}$  in the noncompact case. Note the surprising difference of the value of  $q$  in the case  $C^1$  and  $C^s$ ,  $s \geq 3$ , respectively. Nash's paper (1956) is fundamental not only for the problem under consideration, but also because of the applied method, which gives the foundation of the so-called hard implicit function theorems or Nash-Moser technique. The latter plays an important role in the modern theory of nonlinear partial differential equations. There are many papers concerning this technique, see e.g. Moser (1961, 1966), Jacobowitz (1972), Zehnder (1975), Hamilton (1982) and Hörmander (1985, 1988). Finally we mention the book "Partial Differential Relations" by Gromov (1986), which contains many material and references belonging to our problem. Gromov gives  $q = (n+2)(n+3)/2$  as best value in the smooth case.

One of the main steps in Nash's demonstration (1956) is the solution of the perturbation problem associated to (1), i.e. the determination of  $v : M \rightarrow \mathbb{R}^q$  such that  $F(u+v) = g+f$ , if a solution  $u$  of  $F(u) = g$  is known and  $f$  small in some sense. The linearized equations  $F'(u)v = h$  of our nonlinear problem (1) are  $2du \cdot dv = h$  or locally written

$$\partial_i u \cdot \partial_j v + \partial_j u \cdot \partial_i v = h_{ij}.$$

Here the embedding  $u : M \rightarrow \mathbb{R}^q$  and a symmetric covariant 2-tensor field  $h$  on  $M$  is given and  $v : M \rightarrow \mathbb{R}^q$  is unknown. This is a linear first order system of partial differential equations for  $v$ . The system seems to be very simple, but it does not settle down in a standard class. Now the idea is as follows. Additionally we demand  $\partial_i u \cdot v = 0$ . Then we find the relations

$$\partial_i u \cdot v = 0, \quad \partial_i \partial_j u \cdot v = -h_{ij}/2. \quad (2)$$

which form an algebraic system of  $n(n+3)/2$  equations for  $v$ . Following Gromov and Rohlin (1970) a  $C^2$ -mapping  $u : M \rightarrow \mathbb{R}^q$  is called *free*, if for each  $x \in M$  the  $n(n+3)/2$  vectors  $\partial_i u(x), \partial_i \partial_j u(x)$  of  $\mathbb{R}^q$  are linearly independent. Now let  $u$  be free, then we have a unique solution  $v$  of (2) with minimal pointwise  $\mathbb{R}^q$ -norm, which we will denote by

$$v = -E(u)(0, h)/2.$$

We have the following mapping properties

$$F : C^s \rightarrow C^{s-1}, \quad E(u) : C^s \rightarrow C^s \quad \text{if } u \in C^{s+2}.$$

This shows the loss of differentiability; therefore we cannot use simple successive approximations in order to solve the nonlinear equations. One way out is given by the above mentioned method of hard implicit function theorems consisting in a complicated combination of successive approximations and smoothing operators.

## 2. A New Method to Solve the Perturbation Problem

In this section we give another and short way to handle the perturbation problem, which gives moreover somewhat better results. Let  $M$  be compact. We use certain

Hölder spaces  $C^{s,\lambda}$  of functions or sections of vector bundles over  $M$  with non-negative integers  $s$  and Hölder exponent  $\lambda$ ,  $0 < \lambda < 1$ ; the latter is fixed once for all.  $S_*^{(2)}M$  denotes the bundle of symmetric covariant 2-tensors on  $M$ . Let  $\|E(u)\|_{2,\lambda}$  be the norm of the linear mapping

$$E(u) : C^{2,\lambda}(M, T_* M) \times C^{2,\lambda}(M, S_*^{(2)}M) \rightarrow C^{2,\lambda}(M, \mathbb{R}^q),$$

if  $u$  is a free and smooth mapping. Now we can formulate the following

**Theorem 1.** *Let  $u \in C^\infty(M, \mathbb{R}^q)$  be a free mapping and  $f \in C^{s,\lambda}(M, S_*^{(2)}M)$  with  $s \geq 2$  or  $s = \infty$ . There is a positive number  $\theta$  (independent of  $u, s$  and  $f$ ) with the property: If*

$$\|E(u)\|_{2,\lambda} \|E(u)(0, f)\|_{2,\lambda} \leq \theta, \quad (3)$$

*then there exists a  $v \in C^{s,\lambda}(M, \mathbb{R}^q)$  such that one has*

$$d(u + v) \cdot d(u + v) = du \cdot du + f \quad \text{on } M.$$

**Remark.** Our solubility condition (3) contains two factors coming from the linearized problem taken at the given initial mapping  $u$ . The first depends only on the coefficients of the unknowns of the linearized problem, the second depends only on the solution of the linearized problem with the given perturbation term  $f$ . These two influences must keep a certain balance.

We only sketch the *proof* here. In order to explain the idea, let  $M$  be the  $n$ -dimensional torus. For this special case see also a recent paper of Hörmander (1988), where a much more complicated technique is used. We write our nonlinear equations as

$$\partial_i u \cdot \partial_j v + \partial_j u \cdot \partial_i v + \partial_i v \cdot \partial_j v - f_{ij} = 0. \quad (4)$$

It is well known, that

$$(\mathcal{A} - 1) : C^{s+2,\lambda}(M) \rightarrow C^{s,\lambda}(M)$$

is an isomorphism. Therefore we can apply the operator  $(\mathcal{A} - 1)$  to (4) and obtain after some rearrangement the crucial relations

$$\begin{aligned} & \partial_i \{ (\mathcal{A} - 1)(\partial_j u \cdot v) + \mathcal{A} v \cdot \partial_j v \} + \partial_j \{ (\mathcal{A} - 1)(\partial_i u \cdot v) + \mathcal{A} v \cdot \partial_i v \} \\ & - 2 \{ (\mathcal{A} - 1)(\partial_i \partial_j u \cdot v) + \frac{1}{2} \partial_i v \cdot \partial_j v - \partial_i \partial_j v \cdot \partial_i \partial_j v \\ & + \mathcal{A} v \cdot \partial_i \partial_j v + \frac{1}{2} (\mathcal{A} - 1) f_{ij} \} = 0. \end{aligned} \quad (5)$$

Hence it suffices to solve the new system

$$\partial_i u \cdot v = -(\mathcal{A} - 1)^{-1}(\mathcal{A} v \cdot \partial_i v),$$

$$\partial_i \partial_j u \cdot v = -\frac{1}{2} f_{ij} + (\mathcal{A} - 1)^{-1}(\partial_i \partial_j v \cdot \partial_i \partial_j v - \mathcal{A} v \cdot \partial_i \partial_j v - \frac{1}{2} \partial_i v \cdot \partial_j v).$$

The properties of  $(\mathcal{A} - 1)^{-1}$  guarantee that there is no loss of differentiability, hence

we can solve the system by simple successive approximations if  $f$  fulfils the smallness condition (3). The only difficulty is to treat the  $C^\infty$ -case. We must take care of the constants arising in the inequalities and their dependence on the number  $s$ .

Now some words to the general case. We equip  $M$  with an auxiliary Riemannian metric  $g_0$  of class  $C^\infty$ . The covariant derivatives, the Laplacians as well as the occurring curvature tensors are meant always with respect to  $g_0$ . We have to use the Lichnerowicz-Laplacians for vector fields and symmetric 2-tensor fields

$$\begin{aligned}\Delta_{(1)} t_i &= \Delta t_i - R_{i\cdot}^l t_l, \quad \Delta = \nabla^l \nabla_l, \\ \Delta_{(2)} t_{ij} &= \Delta t_{ij} - 2R_{i\cdot j\cdot}^k t_{kl} - R_{i\cdot}^l t_{lj} - R_{j\cdot}^l t_{il}.\end{aligned}$$

If  $\alpha$  is a sufficiently large real number, then the mappings

$$\begin{aligned}(\Delta_{(1)} - \alpha) : C^{s+2,\lambda}(M, T_* M) &\rightarrow C^{s,\lambda}(M, T_* M), \\ (\Delta_{(2)} - \alpha) : C^{s+2,\lambda}(M, S_*^{(2)} M) &\rightarrow C^{s,\lambda}(M, S_*^{(2)} M)\end{aligned}$$

are isomorphisms. We define a vector field  $N(v)$  and a symmetric 2-tensor field  $L(v)$  by the formulas

$$\begin{aligned}N_i(v) &= -\Delta v \cdot \nabla_i v, \\ L_{ij}(v) &= \nabla^l \nabla_i v \cdot \nabla_l \nabla_j v - \Delta v \cdot \nabla_i \nabla_j v - R_{i\cdot j\cdot}^k \nabla_k v \cdot \nabla_l v - \frac{\alpha}{2} \nabla_i v \cdot \nabla_j v.\end{aligned}$$

One has instead of (5) the identity

$$\begin{aligned}(\Delta_{(2)} - \alpha) \{ \nabla_i u \cdot \nabla_j v + \nabla_j u \cdot \nabla_i v + \nabla_i v \cdot \nabla_j v \} \\ = \nabla_i \{ (\Delta_{(1)} - \alpha)(\nabla_j u \cdot v) - N_j(v) \} + \nabla_j \{ (\Delta_{(1)} - \alpha)(\nabla_i u \cdot v) - N_i(v) \} \\ - 2 \{ (\Delta_{(2)} - \alpha)(\nabla_i \nabla_j u \cdot v) - L_{ij}(v) - (\nabla_i R_{j\cdot}^l + \nabla_j R_{i\cdot}^l - \nabla^l R_{ij})(\nabla_l u \cdot v) \}.\end{aligned}$$

These equations are shown by a straightforward calculation with repeated use of the Ricci identity. This procedure leads to the final system

$$\begin{aligned}\nabla_i u \cdot v &= ((\Delta_{(1)} - \alpha)^{-1} N(v))_i, \\ \nabla_i \nabla_j u \cdot v &= -\frac{1}{2} f_{ij} + ((\Delta_{(2)} - \alpha)^{-1} M(v))_{ij}\end{aligned}\tag{6}$$

with

$$M_{ij}(v) = L_{ij}(v) + (\nabla_i R_{j\cdot}^l + \nabla_j R_{i\cdot}^l - \nabla^l R_{ij})((\Delta_{(1)} - \alpha)^{-1} N(v))_l.$$

Again we can solve (6) with the help of simple successive approximations. For details see Günther (1989a).  $\square$

We give another variant of the perturbation result, which is important for applications in Sect. 3.

**Theorem 2.** *Let  $B \subseteq \mathbb{R}^n$  be the open unit ball and  $B_1, B_2 \subseteq \mathbb{R}^n$  open sets with  $\bar{B}_1 \subseteq B_2$ ,  $\bar{B}_2 \subseteq B$ . Further, let  $u \in C^\infty(\bar{B}, \mathbb{R}^q)$  be a free mapping and  $f = (f_{ij}) \in C^{s,\lambda}(B, \mathbb{R}^{n(n+1)/2})$  with  $s \geq 2$  or  $s = \infty$ . There exists a positive number  $\theta$  (independent of  $u$ ,  $s$  and  $f$ ) with*

the property: If

$$\text{supp } f \subseteq B_1 \quad \text{and} \quad \|E(u)\|_{2,\lambda} \|E(u)(0, f)\|_{2,\lambda} \leq \theta,$$

then there exists a  $v \in C^{s,\lambda}(M, \mathbb{R}^q)$  with

$$\text{supp } v \subseteq B_2 \quad \text{and} \quad \partial_i(u + v) \cdot \partial_j(u + v) = \partial_i u \cdot \partial_j u + f_{ij} \quad \text{in } B.$$

We indicate the *proof* of Theorem 2. We choose a cut-off function  $a \in C^\infty(B)$  with  $a(x) = 1$  for all  $x \in B_1$  and  $a(x) = 0$  for all  $x \in \bar{B} \setminus B_2$ . Further, the Dirichlet problem

$$\Delta w = h \quad \text{in } \bar{B}, \quad w = 0 \quad \text{on } \partial B$$

possesses a unique solution  $w =: \Delta^{-1}h \in C^{s+2,\lambda}(\bar{B})$  if  $h \in C^{s,\lambda}(\bar{B})$ . The definition of  $N_i(v)$ ,  $M_{ij}(v)$  is now modified as follows

$$N_i(v) = 2\partial_i a(\Delta v \cdot v) + a(\Delta v \cdot \partial_i v),$$

$$\begin{aligned} M_{ij}(v) = & \Delta(4\partial_i a \partial_j a(v \cdot v) + 2a\partial_i a(\partial_j v \cdot v) + 2a\partial_j a(\partial_i v \cdot v) + a^2(\partial_i v \cdot \partial_j v)) \\ & - \Delta(a\partial_i \Delta^{-1} N_j(v) + a\partial_j \Delta^{-1} N_i(v) + 3\partial_i a \Delta^{-1} N_j(v) + 3\partial_j a \Delta^{-1} N_i(v)). \end{aligned}$$

Then  $M_{ij}(v)$  is a linear combination of terms

$$\partial^{\alpha_1} a \partial^{\alpha_2} \Delta^{-1} N_i(v) \quad \text{with} \quad |\alpha_1| + |\alpha_2| = 3, |\alpha_2| \leq 2,$$

$$\partial^{\alpha_1} a \partial^{\alpha_2} a(\partial^{\alpha_3} v \cdot \partial^{\alpha_4} v) \quad \text{with} \quad |\alpha_1| + \dots + |\alpha_4| = 4, |\alpha_3|, |\alpha_4| \leq 2$$

and therefore we have  $N_i(v)$ ,  $M_{ij}(v) \in C^{s,\lambda}$  if  $v \in C^{s+2,\lambda}$ . Let  $v$  be a solution of the system

$$\partial_i u \cdot v = -a\Delta^{-1} N_i(v),$$

$$\partial_i \partial_j u \cdot v = \frac{1}{2}(-f_{ij} + \Delta^{-1} M_{ij}(v)) \tag{7}$$

then after some easy calculations it follows that

$$\partial_i(u + a^2 v) \cdot \partial_j(u + a^2 v) = \partial_i u \cdot \partial_j u + a^2 f_{ij} \quad \text{in } B.$$

The properties of  $N_i(v)$ ,  $M_{ij}(v)$  and  $\Delta^{-1}$  guarantee that (7) can be solved by simple approximations too.

### 3. A Result Concerning the Value of $q$

Now let  $M$  be any  $C^\infty$ -manifold, not necessarily compact. We have the following

**Theorem 3.** Let  $g$  be of class  $C^\infty$ , let  $u_0 \in C^\infty(M, \mathbb{R}^q)$  be a free embedding with  $q \geq n(n+3)/2 + 5$  and  $du_0 \cdot du_0 < g$ . Further, let  $\delta$  be any positive continuous function on  $M$ . Then there exists an embedding  $u \in C^\infty(M, \mathbb{R}^q)$  such that

$$du \cdot du = g \quad \text{and} \quad |u(x) - u_0(x)| \leq \delta(x) \quad \text{for every } x \in M.$$

Concerning the existence of free embeddings  $u_0$  we have the following proposition. Its proof is based on the well known theorem of Sard, see e.g. Gromov

and Rohlin (1970, Sect. 2.5). The condition  $du_0 \cdot du_0 < g$  can be reached by easy manipulations.

**Proposition.** *If  $g$  is a continuous metric on  $M$ , then there exists a free embedding  $u_0 \in C^\infty(M, \mathbb{R}^q)$  with  $q = n(n + 5)/2$  and  $du_0 \cdot du_0 < g$  on  $M$ .*

Combining Theorem 3 and the Proposition, we obtain the

**Corollary.** *Every smooth Riemannian manifold possesses a smooth isometric embedding into  $\mathbb{R}^q$  with*

$$q = \max\{n(n + 5)/2, n(n + 3)/2 + 5\}.$$

**Remark.** In general, we have for our  $q$ , that  $q = q_G - 3$  and in special cases, i.e.  $M = S^n$  ( $n$ -dimensional sphere), that  $q = q_G - (n - 2)$ . Here  $q_G$  means that value of  $q$ , which was given by Gromov (1986, Sect. 3.1.7). The main restriction for our  $q$  comes from the existence of a free embedding and not from our method to construct the isometric embedding.

We will give a brief sketch for the proof of Theorem 3. Firstly, one can write

$$g = du_0 \cdot du_0 + \sum_{l \geq 1} h^{(l)}.$$

Thereby  $h^{(l)}$  are symmetric covariant 2-tensor fields on  $M$ , such that in suitable local coordinate systems  $x : u^{(l)} \rightarrow \mathbb{R}^n$

$$h_{11}^{(l)}(x) = a^4(x), \quad h_{ij}^{(l)}(x) = 0 \quad (1 \leq i \leq j \leq n, j \neq 1) \quad (8)$$

with  $a \in C_0^\infty(u^{(l)})$ ; hence  $\text{supp } h^{(l)} \subseteq U^{(l)}$ . Only the one-one-coordinate is different from zero! The family  $\{U^{(l)}\}$  are to be a locally finite covering for  $M$ . Hence it suffices to show:  $u_0$  can be changed into a free embedding  $\tilde{u}_0$  with

$$\begin{aligned} d\tilde{u}_0 \cdot d\tilde{u}_0 &= du_0 \cdot du_0 + h^{(1)} \quad \text{in } U^{(1)}, \\ \tilde{u}_0 &= u_0 \quad \text{in } M \setminus U^{(1)}, \\ |\tilde{u}_0(x) - u_0(x)| &< \delta(x) \quad \text{for all } x \in M. \end{aligned}$$

The construction of such a  $\tilde{u}_0$  will be done in the steps two and three.

Secondly, identifying  $U^{(1)}$  with the open unit ball  $B$  of  $\mathbb{R}^n$ , we choose to every integer  $k \geq 2$  and  $\varepsilon \in (0, 1)$  a free mapping  $u_{\varepsilon, k} \in C^\infty(\bar{B}, \mathbb{R}^q)$  with

$$\begin{aligned} du_{\varepsilon, k} \cdot du_{\varepsilon, k} &= du_0 \cdot du_0 + h^{(1)} + O(\varepsilon^{k+1}), \\ u_{\varepsilon, k} &= u_0 \quad \text{in } B \setminus B_1. \end{aligned}$$

$B_1 \subseteq \mathbb{R}^n$  is an open set with  $\bar{B}_1 \subseteq B$ . The construction of  $u_{\varepsilon, k}$  is made in form of a series

$$u_{\varepsilon, k}(x) = u_0(x) + \varepsilon u_1(\varepsilon, x) + \cdots + \varepsilon^k u_k(\varepsilon, x).$$

(9) means equality up to powers  $\varepsilon^k$ . The proof offers certain technical difficulties; one must take advantage of the special form of  $h^{(1)}$  in (8) as well as the assumption

$q \geq n(n+3)/2 + 5$ . We note, that  $u_{\varepsilon,k}$  has in some sense a singular behaviour as  $\varepsilon \rightarrow 0$ . For details see Günther (1989b, Sects. 3, 4).

Thirdly we apply Theorem 2 with the initial mapping  $u = u_{\varepsilon,k}$ ; the tensor field  $f = f_{\varepsilon,k}$  is determined by the remainder term  $O(\varepsilon^{k+1})$  in (9). From the properties of the  $u_{\varepsilon,k}$  we obtain

$$\|E(u_{\varepsilon,k})\|_{2,\lambda} = O(\varepsilon^{-k_0}), \quad \|f_{\varepsilon,k}\|_{2,\lambda} = O(\varepsilon^{k-2})$$

with an integer  $k_0$  independent of  $k$ . If we choose  $k$  sufficiently large and  $\varepsilon$  sufficiently small, then we can satisfy the solubility condition of Theorem 2. This finishes the proof of Theorem 3.

## References

- Burstin, C. (1931): Ein Beitrag zum Problem der Einbettung der Riemannschen Räume in euklidischen Räumen. Mat. Sb. USSR **38**, 74–85
- Cartan, E. (1927): Sur la possibilité de plonger un espace riemannien donné dans un espace euclidien. Ann. Soc. Polon. Math. **6**, 1–7
- Gromov, M.L., Rohlin, V.A. (1970): Embeddings and immersions in Riemannian geometry. Usp. Mat. Nauk **25**, 3–62 (Russian). [English transl.: Russ. Math. Surv. **25** (1970) 1–57]
- Gromov, M.L. (1986): Partial differential relations. Springer, New York Berlin Heidelberg
- Günther, M. (1989a): On the perturbation problem associated to isometric embeddings of Riemannian manifolds. Ann. Global Anal. Geom. **7**, 69–77
- Günther, M. (1989b): Zum Einbettungssatz von J. Nash. Math. Nachr. **144**, 165–187
- Hamilton, R.S. (1982): The inverse function theorem of Nash and Moser. Bull. Amer. Math. Soc. **7**, 65–222
- Hörmander, L. (1985): On the Nash-Moser implicit function theorem. Ann. Acad. Sci. Fenn. **10**, 255–259
- Hörmander, L. (1988): The Nash-Moser theorem and para-differential operators. Preprint
- Jacobowitz, H. (1972): Implicit function theorems and isometric embeddings. Ann. Math. **95**, 191–225
- Janet, M. (1926): Sur la possibilité de plonger un espace riemannien donné dans un espace euclidien. Ann. Soc. Polon. Math. **5**, 38–43
- Kuiper, N. (1955): On  $C^1$  isometric imbeddings I, II. Proc. Kon. Acad. Wet. Amsterdam A **58**, Indagationes Mathematicae **17**, 545–556, 683–689
- Moser, J. (1961): A new technique for the construction of solutions of non-linear differential equations. Proc. Nat. Acad. Sci. USA **47**, 1824–1831
- Moser, J. (1966): A rapidly convergent iteration method and non-linear partial differential equations I, II. Ann. Scuola Norm. Sup. Pisa **20**, 265–315, 499–535
- Nash, J. (1954):  $C^1$  isometric imbeddings. Ann. Math. **60**, 545–556
- Nash, J. (1956): The imbedding problem for Riemannian manifolds. Ann. Math. **63**, 20–63
- Nash, J. (1966): Analyticity of the solutions of implicit function problems with analytic data. Ann. Math. **84**, 345–355
- Zehnder, E. (1975): Generalized implicit function theorems with applications to some small divisor problems I. Comm. Pure Appl. Math. **28**, 91–140



# On Scattering by Obstacles

Mitsuru Ikawa

Department of Mathematics, Faculty of Science, Osaka University  
Toyonaka, Osaka 560, Japan

## 1. Introduction

Let  $n$  be an odd integer  $\geq 3$ , and let  $\Omega$  be an open bounded set in  $\mathbb{R}^n$  with smooth boundary  $\Gamma$ . We assume that

$$\Omega = \mathbb{R}^n - \overline{\Omega} \quad \text{is connected.}$$

Consider the following acoustic problem

$$\begin{cases} \square u = \frac{\partial^2 u}{\partial t^2} - \Delta u = 0 & \text{in } \Omega \times (-\infty, \infty) \\ u = 0 & \text{on } \Gamma \times (-\infty, \infty) \\ u(x, 0) = f_1(x), \quad \frac{\partial u}{\partial t}(x, 0) = f_2(x). \end{cases} \quad (1.1)$$

It is known that every solution  $u(x, t)$  of (1.1) with finite energy approaches to a solution  $u_0^+(x, t)(u_0^-(x, t))$  of the wave equation in the free space as  $t \rightarrow \infty$  ( $t \rightarrow -\infty$ ). The mapping from the initial data of  $u_0^-(x, t)$  to those of  $u_0^+(x, t)$  is called scattering operator, and the scattering matrix  $\mathcal{S}(\sigma)$  ( $\sigma \in \mathbb{R}$ ) is a representation of this mapping (for the definition, see for example Lax-Phillips [10, p. 170]).  $\mathcal{S}(\sigma)$  is a unitary operator in  $L^2(S^{n-1})$  for all  $\sigma \in \mathbb{R}$  where  $S^{n-1} = \{\omega \in \mathbb{R}^n; |\omega| = 1\}$ , and is of the form

$$\mathcal{S}(\sigma) = I + \mathcal{K}(\sigma) \quad (1.2)$$

where  $I$  denotes the identity operator and  $\mathcal{K}(\sigma)$  is an integral operator. The kernel  $K(\omega, \theta; \sigma)$  of  $\mathcal{K}(\sigma)$  is given in the following way:

Let  $v_-(x, \omega, \sigma)$  be the solution of

$$\begin{aligned} \Delta v + \sigma^2 v &= 0 && \text{in } \Omega, \\ v &= -\exp(-i\sigma\omega) && \text{on } \Gamma, \\ &\cdot v \text{ satisfies the incoming radiation condition.} \end{aligned}$$

Then  $v_-$  has an asymptotic expansion

$$v_-(r\theta, \omega, \sigma) \sim \frac{e^{i\sigma r}}{r^{(n-1)/2}} s(\theta, \omega; \sigma), \quad \text{as } r \rightarrow \infty.$$

Then the kernel  $K(\omega, \theta; \sigma)$  is given by

$$K(\omega, \theta; \sigma) = \left( \frac{\sigma}{2\pi i} \right)^{(n-1)/2} \overline{s(-\theta, \omega; \sigma)}.$$

Roughly speaking, the kernel  $K(\omega, \theta; \sigma)$  represents the rate of the wave reflected by  $\mathcal{O}$  in the direction  $\theta$  for the incident plane wave of frequency  $\sigma$  propagating in the direction  $\omega$ .

Concerning  $\mathcal{S}(\sigma)$ , the following fact is known:

$\mathcal{S}(\sigma)$  is the restriction to the real axis of an  $\mathcal{L}(L^2(S^{n-1}))$ -valued function  $\mathcal{S}(z)$  that is analytic for  $\text{Im } z \leq 0$  and meromorphic in the whole complex plane  $\mathbb{C}$  ([10, p. 166]).

It is an intrinsic subject of scattering theory to consider relationships between the geometry of the obstacle and the analytic property of the scattering matrix. Concerning this subject, the following theorem is fundamental:

**Theorem 5.6 of Chapter V of [10].** *The scattering matrix uniquely determines the scatterer.*

The above theorem shows that all the geometric informations of the obstacle are contained in the scattering matrix. Indeed, all the informations of  $\mathcal{O}$  are in  $\mathcal{S}(z)$ , but how do we extract the geometric informations from the scattering matrix? This is one of the most important and interesting problems of scattering theory.

In this note we would like to consider scattering matrices in connection with the above problem. In Section 2 we present several results on scattering matrices. In Section 3 we propose a conjecture which was given originally by Lax and Phillips in [10]. In Section 4 the validity of this conjecture will be considered for obstacles consisting of several strictly convex bodies.

## 2. Several Results on Scattering Matrices

As to the geometry of obstacles, we introduce the following notion.

**Definition 1.** Suppose that  $\mathcal{O}$  is contained in  $\{x; |x| < \varrho\}$ . We say that  $\mathcal{O}$  is nontrapping if there exists  $T > 0$  such that every broken ray in  $\Omega$  according to the geometric optics starting from a point in  $\Omega(\varrho) = \Omega \cap \{x; |x| < \varrho\}$  necessarily goes out from  $\Omega(\varrho)$  within time  $T$ . We say that  $\mathcal{O}$  is trapping if, for any  $T > 0$ , there is a broken ray staying in  $\Omega(\varrho)$  for a period longer than  $T$ .

## 2.1 Asymptotic Behavior of the Scattering Phase

Since  $\mathcal{S}(\sigma)$  is of trace class,  $\det \mathcal{S}(\sigma)$  is well defined, and the unitarity of  $\mathcal{S}(\sigma)$  implies that  $|\det \mathcal{S}(\sigma)| = 1$ . The scattering phase  $s(\sigma)$  is defined by

$$s(\sigma) = -i \log \det \mathcal{S}(\sigma) \quad \text{for } \sigma \in \mathbb{R}. \quad (2.1)$$

Melrose [13] proved that

$$s(\sigma) = c_n \text{Vol}(\mathcal{O}) \sigma^n + O(\sigma^{n-1}) \quad \text{as } \sigma \rightarrow \pm\infty. \quad (2.2)$$

Analogous to the Weyl formula which is an asymptotics of the distribution of eigenvalues of the interior problem, the volume of the obstacle can be got from the asymptotic behavior of the scattering phase. For nontrapping obstacles Petkov-Popov [17] obtained the full asymptotic expansion of  $s(\sigma)$ .

## 2.2 On the Poles of the Scattering Matrix

As mentioned in the Introduction, the scattering matrix is meromorphic in the whole complex plane. Thus, it is very natural to pose the question:

*How does the geometry of  $\mathcal{O}$  relate to the distribution of poles of  $\mathcal{S}(z)$ ?*

**Purely Imaginary Poles.** Lax and Phillips showed in [11] that there are an infinite number of purely imaginary poles.

**For Nontrapping Obstacles.** By combining the general result in [12] with the results on the propagation of singularities for the problem (1.1) due to [15], we have that, if  $\mathcal{O}$  is nontrapping, there are positive constants  $a$  and  $b$  such that

$$\{z; \text{Im } z \leq a \log(|z| + 1) + b\} \text{ is free from poles of } \mathcal{S}(z). \quad (2.3)$$

**Obstacles Consisting of Two Convex Bodies.** (i) The existence of non-purely imaginary poles was proved for the first time by Bardos-Guillot-Ralston [1]. They considered the case that

$$\mathcal{O} = \mathcal{O}_1 \cup \mathcal{O}_2, \quad \overline{\mathcal{O}_1} \cap \overline{\mathcal{O}_2} = \phi, \quad (2.4)$$

$$\mathcal{O}_1 \text{ and } \mathcal{O}_2 \text{ are strictly convex,} \quad (2.5)$$

and showed that, for any  $\varepsilon > 0$ , the logarithmic domain  $\{z; \text{Im } z \leq \varepsilon \log |z|\}$  contains an infinite number of poles of  $\mathcal{S}(z)$ .

(ii) This result was improved by Ikawa [4] as follows: Let  $a_j \in \Gamma_j = \partial \mathcal{O}_j$ ,  $j = 1, 2$ , be the points such that  $|a_1 - a_2| = \text{dis}(\mathcal{O}_1, \mathcal{O}_2) = d$ . Then,  $\mathcal{S}(z)$  has poles at approximately the points

$$\frac{\pi k}{d} + ic, \quad k \in \mathbb{Z} \quad (2.6)$$

where  $c$  is a positive constant determined by  $d$  and the curvatures of  $\Gamma_j$  at  $a_j$ ,  $j = 1, 2$ .

Later, Gérard [3] gave more precise descriptions of poles.

(iii) Ikawa [6] considered  $\mathcal{O} \subset \mathbb{R}^3$  of the form (2.4) such that the principal curvatures  $\kappa_{j,k}(x)$ ,  $j = 1, 2$ , of  $\Gamma_j$  at  $x$  vanish only at  $a_j$  of order  $2l$ ,  $l \geq 1$ , that is, for  $j = 1, 2$

$$C^{-1}|x - a_j|^{2l} \leq \kappa_{j,k}(x) \leq C|x - a_j|^{2l} \quad \text{for } x \in \Gamma_j, k = 1, 2.$$

In this case there is a sequence of poles  $\{z_j\}_{j=1}^\infty$  such that

$$\operatorname{Im} z_j \rightarrow 0 \quad \text{as } j \rightarrow \infty.$$

**Poles for Convex Obstacles.** For nontrapping obstacles, as (2.3) shows, all the poles of the scattering matrix are over a logarithmic curve. Bardos-Lebeau-Rauch [2] considered strictly convex obstacles with analytic boundary, and showed that, under a certain additional condition, there is a positive constant  $c > 0$  such that  $\{z; \operatorname{Im} z \leq c|z|^{1/3}\}$  contains only a finite number of poles and  $\{z; \operatorname{Im} z \leq (c + \varepsilon)|z|^{1/3}\}$  contains an infinite number of poles for any  $\varepsilon > 0$ .

**Upper Bound of Distribution of Poles.** Melrose showed in [13] the following estimate:

$$\#\{z; \text{poles of } \mathcal{S}(z) \text{ such that } |z| \leq \lambda\} \leq C\lambda^n \quad \text{for all } \lambda > 0. \quad (2.7)$$

*Remark.* Zworski considered in [18, 19] the same problem for scattering by potentials, and showed that (2.7) holds also in this case. Moreover, he showed the estimate (2.7) is optimal in general.

### 3. Modified Lax-Phillips Conjecture

Consider the boundary value problem with parameter  $\mu \in \mathbb{C}$

$$\begin{cases} (\Delta - \mu^2)w(x) = 0 & \text{in } \Omega \\ w(x) = g(x) & \text{on } \Gamma \end{cases} \quad (3.1)$$

for  $g(x) \in C^\infty(\Gamma)$ . It is well known that for  $\operatorname{Re} \mu > 0$  (3.1) has a unique solution  $w(x)$  in  $H^2(\Omega)$ . Denote by  $R(\mu)$  the operator defined by

$$w(x) = (R(\mu)g(\cdot))(x).$$

Then  $R(\mu)$  is an  $\mathcal{L}(L^2(\Gamma), L^2(\Omega))$ -valued holomorphic function in  $\operatorname{Re} \mu > 0$ . By the regularity theorem for elliptic operators,  $R(\mu)$  can be regarded as an operator in  $\mathcal{L}(C^\infty(\Gamma), C^\infty(\overline{\Omega}))$ . As an  $\mathcal{L}(C^\infty(\Gamma), C^\infty(\overline{\Omega}))$ -valued function  $R(\mu)$  can be prolonged analytically to the whole complex plane as a meromorphic function. Concerning the poles of  $\mathcal{S}(z)$  and  $R(\mu)$  we have from Theorem 5.1 of Chapter V of [10]

$$z \text{ is a pole of } \mathcal{S}(z) \text{ if and only if } \mu = iz \text{ is a pole of } R(\mu).$$

Let  $z$  be a pole of  $\mathcal{S}(z)$ . Then  $\mu = iz$  is a pole of  $R(\mu)$ . The above Theorem 5.1 asserts also the existence of a non-trivial  $\mu$ -outgoing solution  $w(x; \mu)$  of (3.1) for

$g \equiv 0$ . Therefore  $u(x, t; \mu) = w(x; \mu) e^{\mu t}$  satisfies (1.1) for  $f_1(x) = w(x; \mu)$ ,  $f_2(x) = \mu w(x; \mu)$ . Thus, if  $\operatorname{Im} z = -\operatorname{Re} \mu$  is very small,  $u(x, t; \mu)$  decays very slowly as  $t \rightarrow \infty$ . This suggests us that, if poles of  $\mathcal{S}(z)$  appear near the real axis, we have probably solutions of (1.1) with finite energy decaying very slowly as  $t \rightarrow \infty$ . Thus, we suspect that the stronger solutions to (1.1) are trapped by obstacle  $\mathcal{O}$ , the nearer to the real axis the poles of  $\mathcal{S}(z)$  will appear, and we suspect also that the stronger rays of geometric optics in  $\Omega$  are trapped, the stronger solutions to (1.1) will be trapped by  $\mathcal{O}$ .

Thus, we would like to propose the following conjecture:

**Modified Lax-Phillips Conjecture.** *If  $\mathcal{O}$  is trapping, there exists  $\alpha > 0$  such that  $\mathcal{S}(z)$  has an infinite number of poles in  $\{z; \operatorname{Im} z \leq \alpha\}$ .*

Recall that, for nontrapping obstacle  $\mathcal{O}$ , all the poles of  $\mathcal{S}(z)$  are over a logarithmic curve  $\operatorname{Im} z = a \log(|z| + 1) + b$ . Then, for any  $\alpha > 0$  there are only a finite number of poles in  $\{z; \operatorname{Im} z \leq \alpha\}$ . Thus, if the above conjecture is true, the existence of such  $\alpha$  becomes a characterization of trapping obstacles by means of distribution of poles of scattering matrices.

Hereafter, we say that MLPC (abbreviation of the Modified Lax-Phillips Conjecture) is valid for obstacle  $\mathcal{O}$  when there is  $\alpha > 0$  such that the scattering matrix  $\mathcal{S}(z)$  corresponding to  $\mathcal{O}$  has an infinite number of poles in  $\{z; \operatorname{Im} z \leq \alpha\}$ .

Even though it is more than 20 years since the original conjecture was given, the examples of obstacles for which the validity of MLPC is proved are few. As presented in Section 2, obstacles consisting of two convex bodies were known as these examples.

#### 4. Obstacles Consisting of Several Strictly Convex Bodies

As mentioned in the previous section, MLPC is valid for  $\mathcal{O}$  consisting of two strictly convex bodies. In this section we consider an extension of this result to obstacles consisting of several strictly convex bodies. Here we would like to mention about the geometrical difference between  $\mathcal{O}$  consisting of two strictly convex bodies and  $\mathcal{O}$  consisting of more than two bodies. For  $\mathcal{O}$  consisting of two strictly convex bodies, there is only one primitive periodic rays in  $\Omega$ . On the other hand, when  $\mathcal{O}$  is consisted of three convex bodies for example, there are infinitely many primitive periodic rays in  $\Omega$  in general. The infiniteness of the number of primitive periodic rays in  $\Omega$  makes the problem difficult to treat. In this case we have to control the complexity caused by the infiniteness of the number of primitive periodic rays in  $\Omega$ . It seems to us that the asymptotic behavior of periodic rays is closely related to the ergodic property of rays in  $\Omega$ . Actually we can control the complexity of periodic rays only for obstacles consisting of several small balls.

Now we shall state our theorem. Let  $P_j$ ,  $j = 1, 2, \dots, L$  ( $L \geq 3$ ), be points in  $\mathbb{R}^3$ . For  $\varepsilon > 0$  we set

$$\mathcal{O}_\varepsilon = \bigcup_{j=1}^L \mathcal{O}_{j,\varepsilon}, \quad \mathcal{O}_{j,\varepsilon} = \{x; |x - P_j| < \varepsilon\}.$$

**Theorem 4.1.** Assume that  $P_j$ ,  $j = 1, 2, \dots, L$  ( $L \geq 3$ ), satisfy the condition

$$\text{any triple of } P_j \text{'s does not lie on a straight line.} \quad (\text{A})$$

Then, there is a positive constant  $\varepsilon_0$  such that, for all  $0 < \varepsilon \leq \varepsilon_0$ , MLPC is valid for  $\mathcal{O}_\varepsilon$ .

The proof of Theorem 4.1 will be devided into several steps.

## 4.2 General Theorem

We present a theorem given in [8]. Let  $\mathcal{O}_j$ ,  $j = 1, 2, \dots, L$  ( $L \geq 3$ ), be bounded open sets  $\mathbb{R}^3$  with smooth boundary  $\Gamma_j$  satisfying

$$\text{every } \mathcal{O}_j \text{ is strictly convex,} \quad (\text{H.1})$$

$$\begin{aligned} \text{for every } \{j_1, j_2, j_3\} \in \{1, 2, \dots, L\}^3 \text{ such that } j_l \neq j_{l'} \text{ if } l \neq l', \\ (\text{convex hull of } \overline{\mathcal{O}_{j_1}} \text{ and } \overline{\mathcal{O}_{j_2}}) \cap \overline{\mathcal{O}_{j_3}} = \phi. \end{aligned} \quad (\text{H.2})$$

We set

$$\mathcal{O} = \bigcup_{j=1}^L \mathcal{O}_j, \quad \Gamma_j = \partial \mathcal{O}_j. \quad (4.1)$$

Denote by  $\gamma$  an oriented periodic ray in  $\Omega = \mathbb{R}^3 - \overline{\mathcal{O}}$ , and we shall use the following notations:

- $d_\gamma$ : the length of  $\gamma$ ,
- $T_\gamma$ : the primitive period of  $\gamma$ ,
- $i_\gamma$ : the number of the reflecting points of  $\gamma$ ,
- $P_\gamma$ : the linearized Poincaré map of  $\gamma$ .

We define a function  $F_D(s)$  ( $s \in \mathbb{C}$ ) by

$$F_D(s) = \sum_{\gamma} \frac{(-1)^{i_\gamma} T_\gamma}{|I - P_\gamma|^{1/2}} e^{-sd_\gamma} \quad (4.2)$$

where the summation is taken over all the oriented periodic rays in  $\Omega$  and  $|I - P_\gamma|$  denotes the determinant of  $I - P_\gamma$ .

Concerning the periodic rays in  $\Omega$  we have

$$\#\{\gamma; \text{periodic ray in } \Omega \text{ such that } d_\gamma < r\} < e^{a_0 r} \quad (4.3)$$

and

$$|I - P_\gamma| \geq e^{2a_1 d_\gamma}, \quad (4.4)$$

where  $a_0$  and  $a_1$  are positive constants depending on  $\mathcal{O}$ . The estimates (4.3) and (4.4) imply that the right hand side of (4.2) converges absolutely in  $\{s \in \mathbb{C}; \operatorname{Re} s > a_0 - a_1\}$ . Thus  $F_D(s)$  is well defined in  $\{s \in \mathbb{C}; \operatorname{Re} s > a_0 - a_1\}$ , and holomorphic in this domain.

Now we have

**Theorem 4.2.** Let  $\mathcal{O}$  be an obstacle given by (4.1) satisfying (H.1) and (H.2). If  $F_D(s)$  cannot be prolonged analytically to an entire function, then MLPC is valid for  $\mathcal{O}$ .

The proof is based on the trace formula due to Bardos-Guillot-Ralston [1]. The essential part of the proof is given in [8, Section 2].

### 4.3 Zeta Functions of Symbolic Flows

In order to consider singularities of  $F_D(s)$  we shall use the fact that  $F_D(s)$  has a close relationship to a zeta function of a symbolic flow. Denote by  $v_0$  the abscissa of the convergence of the right hand side of (4.2), that is,

$$v_0 = \inf\{v; \text{ the right hand side of (4.2) converges absolutely for } \operatorname{Re} s > v\}.$$

We shall show the following

**Proposition 4.3.** Let  $\zeta(s)$  be the zeta function defined by (4.6). Then, we have that

$$F_D(s) - \left(-\frac{d}{ds} \log \zeta(s)\right) \text{ is holomorphic in } \operatorname{Re} s > v_0 - a_2$$

where  $a_2$  is a positive constant depending on  $\mathcal{O}$ .

Proposition 4.3 implies that singularities in  $\operatorname{Re} s > v_0 - a_2$  of  $F_D(s)$  coincide with those of  $\zeta(s)$ . Therefore, in order to apply Theorem 4.2 it suffices to show the existence of singularities of  $\zeta(s)$  in  $\operatorname{Re} s > v_0 - a_2$ .

We introduce some notations of symbolic flows. Let  $A = (A(i,j))_{i,j=1,2,\dots,L}$  be the  $L \times L$  matrix defined by  $A(i,j) = 1(i \neq j)$  and  $A(j,j) = 0$ . Following Parry-Pollicott [16] we set

$$\Sigma_A = \{\xi = (\dots, \xi_{-1}, \xi_0, \xi_1, \dots); \xi_j \in \{1, 2, \dots, L\} \text{ and } A(\xi_j, \xi_{j+1}) = 1 \text{ for all } j\}.$$

Denote by  $\sigma_A$  the shift transformation defined by

$$(\sigma_A \xi)_j = \xi_{j+1}.$$

We consider relationships between  $\Sigma_A$  and bounded broken rays in the outside of  $\mathcal{O}$ . Let  $(\dots, l_{-1}, l_0, l_1, \dots)$  be the reflection order of a bounded broken ray in  $\Omega$ . Then, as was shown in [5], it is an element of  $\Sigma_A$ . Conversely, for each element of  $\xi \in \Sigma_A$  there exists a unique broken ray with the reflection order  $\xi$ . Note that a periodic ray in  $\Omega$  corresponds to a periodic element  $\xi \in \Sigma_A$ , that is,  $\sigma_A^n \xi = \xi$  for some  $n$ .

We shall define real valued functions  $f$  and  $g$  on  $\Sigma_A$ . We set

$$f(\xi) = |X_0 X_1|$$

where  $X_j$  denote the  $j$ -th reflection point of the broken ray corresponding to  $\xi$ . Suppose that  $\xi \in \Sigma_A$  satisfies  $\sigma_A^n \xi = \xi$  for some  $n$ . Set  $\mathbf{i} = (\xi_0, \dots, \xi_{n-1})$ , and let

$\varphi_{i,0}^\infty$  be the phase function defined in [5, Section 5]. Denote by  $\lambda_1(\xi)$  and  $\lambda_2(\xi)$  the eigenvalues of  $P_\gamma$  greater than 1, and by  $\kappa_l(\xi), l = 1, 2$ , the principal curvatures at  $X_0$  of the wave front of the phase function  $\varphi_{i,0}^\infty$ . Then we have

$$\lambda_1(\xi)\lambda_2(\xi) = \prod_{j=1}^n (1 + f(\sigma_A^j \xi) \kappa_1(\sigma_A^j \xi))(1 + f(\sigma_A^j \xi) \kappa_2(\sigma_A^j \xi)). \quad (4.5)$$

Define  $g(\xi)$  for an periodic element  $\xi$  by

$$g(\xi) = -\frac{1}{2} \log(1 + f(\xi) \kappa_1(\xi))(1 + f(\xi) \kappa_2(\xi)).$$

By using the fact that the periodic elements are dense in  $\Sigma_A$ , we can extend  $g(\xi)$  for all  $\xi \in \Sigma_A$  by the continuity.

Define  $\zeta(s)$  by

$$\zeta(s) = \exp \left( \sum_{n=1}^{\infty} \frac{1}{n} \sum_{\sigma_A^n \xi = \xi} \exp S_n(-sf(\xi) + g(\xi) + \pi i) \right) \quad (4.6)$$

where we set

$$S_n(-sf(\xi) + g(\xi) + \pi i) = \sum_{k=0}^{n-1} (-sf(\sigma_A^k \xi) + g(\sigma_A^k \xi) + \pi i).$$

It is easy to see that the right hand side of (4.6) converges for  $\operatorname{Re} s$  large. We call  $\zeta(s)$  a zeta function of the symbolic flow  $(\Sigma_A, \sigma_A)$ . Then, we have the following relation for  $\operatorname{Re} s > v_0$ :

$$\begin{aligned} F_D(s) &= \left( -\frac{d}{ds} \log \zeta(s) \right) \\ &= \sum_{\gamma} T_\gamma(-1)^{\ell_\gamma} \{ |I - P_\gamma|^{-1/2} - (\lambda_1 \lambda_2)^{-1/2} \} \exp(-sd_\gamma). \end{aligned} \quad (4.7)$$

Note that the following estimate holds:

$$| |I - P_\gamma|^{-1/2} - (\lambda_1 \lambda_2)^{-1/2} | \leq C(\lambda_1 \lambda_2)^{-1/2} e^{-a_2 d_\gamma} \quad (4.8)$$

where  $a_2$  is a positive constant depending on  $\mathcal{O}$ . Then, we have immediately Proposition 4.3 by substituting (4.8) in the right hand side of (4.7).

#### 4.4 Singular Perturbations of Symbolic Flows

It is likewise difficult in general to show the existence of singularities of the zeta function  $\zeta(s)$ , because there is no  $s \in \mathbb{C}$  such that  $-sf(\xi) + g(\xi) + \pi i$  is real for all  $\xi \in \Sigma_A$ . Namely, this fact makes impossible to apply the Ruelle-Perron-Frobenious theorem to get the existence of poles of  $\zeta(s)$ .

When the bodies composing  $\mathcal{O}$  are small in comparison with the distances between each other,  $\zeta(s)$  can be approximated by a zeta function of a graph, which

is much easier to treat. To this end, we use theorems on singular perturbations of symbolic flows proved in [7, 9]. Denote by  $\zeta_\varepsilon$ ,  $f_\varepsilon$ ,  $g_\varepsilon$  the  $\zeta$ ,  $f$ ,  $g$  attached to  $\mathcal{O}_\varepsilon$  respectively. If we set

$$f_0(\xi) = |P_{\xi_0} P_{\xi_1}|, \quad \tilde{g}_0(\xi) = \frac{1}{2} \log \left( \frac{1}{4} \cos \frac{\Theta(\xi)}{2} \right)$$

where  $\Theta(\xi)$  denotes the angle  $P_{\xi_{-1}} P_{\xi_0} P_{\xi_1}$ , it holds that

$$|\log \varepsilon| ||| f_\varepsilon - f_0 |||_\theta, ||| g_\varepsilon - (\log \varepsilon + \tilde{g}_0) |||_\theta \rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0$$

for every fixed  $\theta > 0$  (For the definition of the norm  $||| \cdot |||_\theta$ , see [16]). By using the above relations we have the following expression of  $\zeta_\varepsilon(s)$ :

$$\zeta_\varepsilon(s) = Z_\varepsilon(s - (\log \varepsilon + \pi i)/d_{\max}).$$

Here  $d_{\max} = \max_{j \neq k} |P_j P_k|$  and  $Z_\varepsilon(\tilde{s})$  is a function defined by

$$Z_\varepsilon(\tilde{s}) = \exp \left( \sum_{n=1}^{\infty} \frac{1}{n} \sum_{\sigma_A^n \xi = \tilde{s}} \exp S_n r_\varepsilon(\xi, \tilde{s}) \right),$$

$$r_\varepsilon(\xi, \tilde{s}) = -\tilde{s} f_\varepsilon(\xi) + h_\varepsilon(\xi) + k(\xi) \log \varepsilon,$$

$$k(\xi) = 1 - f_0(\xi)/d_{\max},$$

where  $h_\varepsilon(\varepsilon \geq 0)$  satisfies

$$||| h_\varepsilon - h_0 |||_\theta \rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0,$$

$$h_0(\xi) \text{ is real if } k(\xi) = 0.$$

Then Theorem 1 of [9] guarantees that  $Z_\varepsilon(\tilde{s})$  has a pole near  $s_0 \in \mathbb{R}$  when  $\varepsilon$  is small.<sup>1</sup> Thus, we have the existence of pole of  $\zeta_\varepsilon(s)$ . Thus Theorem 4.1 is proved.

## References

1. Bardos C., Guillot J.C., Ralston J.: La relation de Poisson pour l'équation des ondes dans un ouvert non borné. Application à la théorie de la diffusion. Comm. Partial Diff. Equations **7** (1982) 905–958
2. Bardos C., Lebeau G., Rauch J.: Scattering frequencies and Gevrey 3 singularities. Invent. math. **90** (1987) 77–114
3. Gérard C.: Asymptotique des poles de la matrice de scattering pour deux obstacles strictement convexes. Bull. Soc. Math. France **116** (31) (1989)
4. Ikawa M.: On the poles of the scattering matrix for two strictly convex obstacles. J. Math. Kyoto Univ. **23** (1983) 127–194
5. Ikawa M.: Decay of solutions of the wave equation in the exterior of several convex bodies. Ann. Inst. Fourier **38** (1988) 113–146
6. Ikawa M.: Trapping obstacles with a sequence of poles of the scattering matrix converging to the real axis. Osaka J. Math. **22** (1985) 657–689

<sup>1</sup> In [7] the existence of poles is proved under some additional assumptions.

7. Ikawa M.: Singular perturbation of symbolic flows and poles of the zeta functions. *Osaka J. Math.* **27** (1990) 281–300
8. Ikawa M.: On the distribution of poles of the scattering matrix for several convex bodies. Proc. of the Conference in Honor of Prof. T. Kato “Functional Analytic Methods for Partial Differential Equations”. Lecture Notes in Mathematics, vol. 1450. Springer, Berlin Heidelberg New York 1990, pp. 210–225
9. Ikawa M.: On the existence of poles of zeta functions for certain symbolic dynamics. In preparation
10. Lax P.D., Phillips R.S.: Scattering theory. Revised edition. Academic Press, New York 1989
11. Lax P.D., Phillips R.S.: Decaying modes for the wave equation in the exterior of an obstacle. *Comm. Pure Appl. Math.* **22** (1969) 737–787
12. Lax P.D., Phillips R.S.: A logarithmic bound on the location of the poles of the scattering operator. *Arch. Rat. Mech.* **40** (1971) 268–280
13. Melrose R.: Polynomial bound on the distribution of poles in scattering by an obstacle. *Journées Équations aux Dérivées Partielles*, St. Jean de Monts, 1984
14. Melrose R.: Weyl asymptotics for the phase in obstacle scattering. *Comm. Partial Diff. Equations* **13** (1988) 1431–1439
15. Melrose R., Sjöstrand J.: Singularities of boundary problems, I and II. *Comm. Pure Appl. Math.* **31** (1978) 593–617; **35** (1982) 129–168
16. Parry W., Pollicott M.: An analogue of the prime number theorem for closed orbits of Axiom A flows. *Ann. Math.* **118** (1983) 573–591
17. Petkov V., Popov G.: Asymptotic behavior of the scattering phase for nontrapping obstacles. *Ann. Inst. Fourier* **32** (1982) 111–149
18. Zworski M.: Sharp polynomial bounds on the number of scattering poles of radial potentials. *J. Funct. Anal.* **82** (1989) 370–403
19. Zworski M.: Sharp polynomial bounds on the number of scattering poles. *Duke Math. J.* **59** (1989) 311–323

# Interaction des Singularités Faibles Pour les Équations d’Ondes Semi-linéaires

Gilles Lebeau

Université de Paris-Sud, Centre d’Orsay, Département de Mathématiques, Bâtiment 425  
F-91405 Orsay Cedex, France

## I. Introduction

Soit  $\square = \partial_t^2 - \Delta_x$  l’opérateur des ondes, où  $t \in \mathbb{R}$ ,  $x \in \mathbb{R}^d$ , et  $\Omega$  un ouvert de  $\mathbb{R}^{1+d}$  qui est un domaine de détermination pour  $\omega = \Omega \cap \{t = 0\}$ .

Soit  $u(t, x)$  une fonction réelle continue sur  $\Omega$  appartenant localement à l’espace

$$C^0(\mathbb{R}_t, H^{s+1}(\mathbb{R}^d)) \cap C^1(\mathbb{R}_t, H^s(\mathbb{R}^d))$$

où  $H^s$  est l’espace de Sobolev usuel et vérifiant dans  $\Omega$  l’équation des ondes semi-linéaires

$$(1) \quad \begin{cases} \square u = F(t, x, u, \nabla u) & \nabla u = (\partial_t u, \nabla_x u) \\ u|_{t=0} = u_0 \in H_{\text{loc}}^{s+1}(\omega) & \frac{\partial u}{\partial t}|_{t=0} = u_1 \in H_{\text{loc}}^s(\omega) \end{cases}$$

où  $F$  est une fonction  $C^\infty$  de ses arguments et où  $s > d/2$  (ou  $s > (d/2) - 1$  si  $F$  ne dépend pas de  $\nabla u$ ).

On s’intéresse à déterminer les singularités de la solution  $u$  de (1) dans  $\Omega_+ = \Omega \cap \{t > 0\}$  en fonction de ses données de Cauchy  $u_0$  et  $u_1$  [ou encore en fonction de  $u|_{\Omega_-}$ ,  $\Omega_- = \Omega \cap \{t < 0\}$ ]. Plus précisément si  $p \in T^*\Omega_+ \setminus \Omega_+$ , on cherche à déterminer pour quelles valeurs de  $\sigma$  a-t-on

$$(2) \quad u \in H_p^\sigma \quad (\text{espace de Sobolev microlocal}).$$

Depuis les travaux de pionnier de J.-M. Bony [8-13] ce type de problème a été intensivement étudié, et son contenu s’est avéré très riche.

Rappelons les résultats généraux sur les solutions du (1).

(3) **Théorème 1** (J.-M. Bony [13]). Soit  $u$  solution de (1) vérifiant  $u \in H_{\text{loc}}^{s_0}(\Omega)$ ,  $s_0 = (1 + d)/2 + 1 + \varrho$ ,  $\varrho > 0$ . Si  $p$  est non caractéristique, on a  $u \in H_p^{s_0 + \varrho + 1}$ . Si  $p$  est caractéristique et si la bicaractéristique de  $\square$  passant par  $p$  est issue de  $q = (x_0, \xi_0) \in T^*\omega \setminus \omega$ , on a  $u \in H_p^\sigma$  pour  $\sigma \leq s_0 + \varrho$  si les données vérifient  $u_j \in H_q^{\sigma-j}$ .

En d’autres termes, il n’y a pas d’effets non linéaires jusqu’à la régularité  $s_0 + \varrho$ .

(4) **Théorème 2** (M. Beals [4], J.-Y. Chemin [17]). *Sous les mêmes hypothèses, si  $\gamma$  est une bicaractéristique de  $\square$ ,  $p_1, p_2$  deux points de  $\gamma$ , on a  $u \in H_{p_1}^\sigma$  ssi  $u \in H_{p_2}^\sigma$  pour  $\sigma \leq 3s_0 - (1+d) - 2 = s_0 + 2\varrho$ .*

Les Théorèmes 1 et 2 ont leur analogue pour les équations totalement non linéaires d'ordre quelconque à linéarisé strictement hyperbolique (voir [8, 17]).

Dans ce cadre, le Théorème 2 exprime que jusqu'à la régularité  $s_0 + 2\varrho$  les effets non linéaires se réduisent à une opération de somme fibre à fibre sur le front d'onde.

Pour les régularités plus élevées, il n'y a pas de théorème de nature purement géométrique sur le front d'onde. Dans [4], M. Beals a construit une solution  $u \in H^s$  de  $\square u + \beta u^3 = 0$  ( $\beta \in C^\infty$ ), dont les données de Cauchy sont  $C^\infty$  hors de l'origine et  $u \notin H^{3s-d+2+\varepsilon}$  à l'intérieur du cône d'onde, avec  $d \geq 2$ . [Le cas de la dimension 1 d'espace a été élucidé par J. Rauch et M. Reed [36], et J.-Y. Chemin [16] qui donnent un résultat complet jusqu'à  $\sigma = \infty$ .]

Il est toutefois possible d'obtenir des résultats de localisation du front d'onde de  $u$  jusqu'à  $\sigma = \infty$  en faisant des hypothèses de conormalité sur les données de Cauchy ou sur la solution  $u$  dans le passé  $\Omega_-$ . Dans ce cadre, le théorème d'interaction de trois ondes progressives a été obtenu simultanément par J.-M. Bony [11] et R. Melrose et N. Ritter [32] et [33], puis généralisé et amélioré par J.-Y. Chemin [18] et Sa-Baretto [41].

La difficulté principale de ce type de problème provient du fait que la géométrie qui porte les singularités de la solution est elle-même singulière. Les techniques utilisées font intervenir soit des microlocalisations d'ordre supérieur (J.-M. Bony) soit des espaces de distributions conormales définis par éclatement (R. Melrose).

On se propose ici de décrire certains résultats qu'on peut obtenir en faisant des hypothèses d'analyticité sur la géométrie qui porte les singularités. Dans cette direction, on a :

(5) **Théorème 3** [26] ( $d = 3$ ). *Soit  $u \in H^s(\Omega)$ ,  $s > 2$  vérifiant (1) où  $F(t, x, u)$  est polynomial en  $u$  et où les données  $u_0, u_1$  sont des distributions intégrales de Fourier  $C^\infty$  sur  $A$ , lagrangienne analytique réelle lisse de  $T^*\omega$ . Pour tout réel  $\sigma$ , il existe un ensemble sous-analytique, homogène isotrope  $L_\sigma$  de  $T^*\Omega$  tel que  $WF(u) \subset L_\sigma$ . En particulier, pour tout  $k$ ,  $u$  est de classe  $C^k$  sur un ouvert dense de  $\Omega$ .*

La preuve du Théorème 3, qui utilise le théorème de désingularisation de Hironaka, ne fournit pas d'estimation sur  $L_\sigma$ . Les résultats qui suivent fournissent pour les solutions de (1) une majoration de  $L_\sigma$  et comme corollaire, la preuve d'une conjecture de J.-M. Bony sur le pincement d'une onde progressive en dimension 2 d'espace. (Démontré dans [27] pour  $F$  polynomiale en  $u$ .)

(6) **Théorème 4** ( $d = 2$ ). *Soit  $u$  solution de (1) avec  $u_0$  et  $u_1$  conormales classiques sur  $A = T_V^*\omega$ , où  $V$  est une courbe analytique lisse possédant un minimum non dégénéré de son rayon de courbure en  $A \in \mathbb{R}^2$ . Soit  $S$  la queue d'aronde issue de  $V$  et  $Q_+$  le demi-cône d'onde d'avenir issu du point de pincement  $B$  de  $S$ . Alors près de  $B$ , on a  $WF(u) \subset T_{S \cup Q_+}^*$ .*

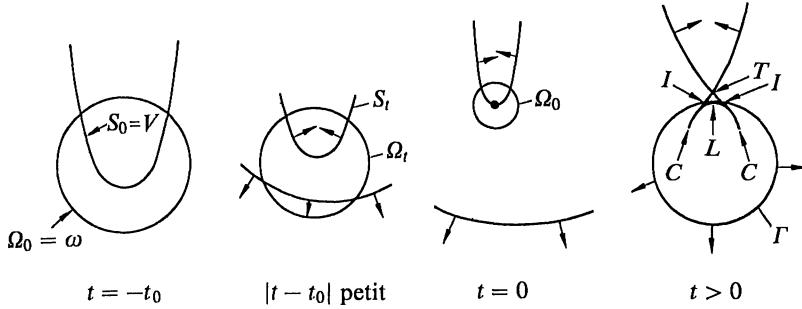


Fig. 1. Pincement d'une onde progressive

Ce résultat a été amélioré par J.-M. Delort [23] qui prouve en particulier que  $u$  est conormale aux points lisses de  $T_{S \cup Q_+}^*$ .

On n'abordera pas ici les travaux relatifs aux équations totalement non linéaires (voir Alinhac [1, 2, 3]), ni les problèmes aux limites ([6, 7, 21, 30, 43, 44, 45, 46]). Signalons à cet égard que R. Melrose, Zworski, Sa-Baretto ont obtenu le théorème de diffraction non-linéaire d'une onde conormale par un obstacle convexe.

## II. Calcul Multilinéaire

La première étape pour obtenir par exemple le Théorème 4 est de majorer le front d'onde de  $u$  par le front d'onde de distributions construites explicitement à partir des données  $u_0, u_1$ , et qui décrivent les phénomènes de propagation-interaction.

On appelle diagramme  $D$  la donnée d'un ensemble fini  $I = \{1, \dots, N\}$  muni d'une partition et d'une application  $f : I \rightarrow I \cup \{0\}$  telle que  $f(I_k) \subset I_{k-1}$ ,  $f(I_1) = \{0\}$ . On définit le degré de  $D$  par  $\deg D = \text{card } J$ ,  $J = \{i \in I, f^{-1}(i) = \phi\}$ . Soit  $e_+(z)$ ,  $(z = (t, x))$  la solution élémentaire de  $\square$  à support dans l'avenir. Pour  $(z_0, z_1, \dots, z_N) \in \mathbb{R}^{(1+d)(N+1)}$  on pose:

$$(8) \quad \begin{cases} [D] = \prod_{i \notin J} \nabla^{\beta_i} e_+(z_{f(i)} - z_i) \prod_{i \in J} e_+(z_{f(i)} - z_i) \\ \{D\} = \prod_{i \in J} (v_i^0(x_i) \delta'_{t_i=0} + v_i^1(x_i) \delta_{t_i=0}). \end{cases}$$

Ici  $\beta_i$  vérifie  $|\beta_i| \leq 1$ ,  $v_i^0 \in \text{vect}\{\mathcal{U}_1, \nabla^\beta \mathcal{U}_0, |\beta| \leq 1\}$  et  $v_i^1 \in \text{vect}\{\nabla^\beta \mathcal{U}_1, \nabla^\gamma \mathcal{U}_0, |\gamma| \leq 1, |\gamma| \leq 2\}$ , et on omet dans la notation la dépendance en  $\beta_i, v_i^0$ . Avec les techniques de [27, 29] on prouve:

(9) **Théorème 5.** Soit  $u$  vérifiant (1) avec  $s = (d/2) + \varrho$ ,  $\varrho > 0$ . Pour  $M \geq 1$  et  $\sigma \in [s + 1 + (M - 1)\varrho, s + 1 + M\varrho[$  on a:

$$(10) \quad \begin{aligned} WF^\sigma(u)|_{t>0} &\subset \\ &\left\{ (z_0, \zeta_0); \text{ il existe un diagramme } D \text{ avec } \deg(D) \leq M \right. \\ &\left. \text{et } (z_0, z_1, \dots, z_N, \zeta_0, 0, \dots, 0) \in WF[[D] \cdot \{D\}] \right\}. \end{aligned}$$

La base de la preuve du Théorème 5 consiste à utiliser une modification du paraproduit de Bony. L'idée est d'utiliser des décompositions d'une fonction  $f(y)$  de la forme

$$(11) \quad f = f_{1,\lambda} + f_{2,\lambda}; \quad f_{1,\lambda} = (2\pi)^{-n} \int_{|\eta| \leq c_0 \lambda^\delta} e^{iy\eta} \widehat{f}(\eta) d\eta$$

où le grand paramètre  $\lambda$  est la fréquence du phénomène qu'on observe. Si  $f \in H^s(\mathbb{R}^n)$ ,  $s = (n/2) + \varrho$ ,  $\varrho > 0$ , on a alors à la fois, avec  $\sigma = s - \nu$ ,  $\nu \in ]0, \varrho[$

$$(12) \quad \|\partial_x^\alpha f_1\|_s \leq C_\alpha \lambda^{|\alpha|\delta} \quad \|f_2\|_\sigma \leq C \lambda^{-\nu\delta}.$$

(Le paraproduit correspondrait à  $\delta = 1$  et  $C_0$  petit.)

Pour une distribution  $f(y, \lambda)$  dépendant d'un paramètre  $\lambda \geq 1$ , telle que  $|\widehat{\varphi f}(\eta, \lambda)| \leq \text{polynôme}(\eta, \lambda)$  pour tout  $\varphi \in C_0^\infty$ , on définit le  $\lambda$ -WF de  $f$  par

$$(y_0, \eta_0) \notin \lambda\text{-WF}(f) \Leftrightarrow \exists \varphi \text{ localisant près de } y_0,$$

$$(13) \quad V \text{ voisinage de } \eta_0 \text{ tels que } \int_{\eta \in \lambda V_0} |\widehat{\varphi f}(\eta, \lambda)|^2 d\eta \in \mathcal{O}(\lambda^{-\infty}).$$

Si  $C_p = \{2^{p-1} \leq |\eta| < 2^{p+1}\}$ , et  $\mu \in \mathbb{R}$  on pose

$$(14) \quad a_p^2(f) = \int_{C_p} |\widehat{f}(\eta, 2^p)|^2 d\eta \quad |f|_\mu^2 = \sum_{p \geq 0} 2^{2p\mu} a_p^2(f).$$

Alors si  $t > n/2$  et  $a(x, \lambda)$  vérifie pour un  $\delta \in ]0, 1[$   $\|\partial_x^\alpha a\|_t \leq C_\alpha \lambda^{|\alpha|\delta}$ , on a  $\lambda\text{-WF}(af) \subset \lambda\text{-WF}(f)$  et pour  $v(x) \in H^\mu$ ,  $|av|_\mu \leq \text{cte} \|v\|_\mu$  pour tout  $\mu$ , et on a:

(15) **Lemme 1.** Supposons donné pour tout  $\lambda$  une décomposition  $g(y) = g_1(y, \lambda) + g_2(y, \lambda)$  avec  $(y_0, t\eta_0) \notin \lambda\text{-WF}(g_1)$ ,  $\forall t > 0$  et  $|g_2|_\mu < \infty$ . Alors  $g \in H_{(y_0, \eta_0)}^\mu$ .

D'où on déduit par exemple par la formule de Taylor en choisissant  $\delta = 1 - 0$  et  $\sigma = (n/2) + 0$  dans (12).

(16) **Proposition 1** [29]. Soit  $u \in H^s(\mathbb{R}^n)$ ,  $s = (n/2) + \varrho$ ,  $\varrho > 0$  et  $F \in C^\infty$ . Pour  $N \geq 1$  et  $\mu \in [(d/2) + N\varrho, d/2 + (N+1)\varrho[$  on a

$$WF^{(\mu)}(F(u)) \subset \bigcup_{j=1}^N WF^\mu(u^j).$$

### III. Deuxième Microlocalisation

La difficulté à pouvoir tirer de l'information à partir du Théorème 5 est que le produit  $[D] \cdot \{D\}$  est caractéristique, c'est-à-dire le front d'onde du produit tensoriel  $[D] \otimes \{D\}$  rencontre le fibré conormal à la diagonale. C'est à ce stade qu'intervient la deuxième microlocalisation. C'est M. Kashiwara qui, dans les années 70, a l'idée

d'introduire les 2-microfonctions associées à une sous-variété involutive du fibré cotangent. Dans ce cadre, le calcul symbolique des opérateurs 2-microdifférentiels a été développé par Y. Laurent [25]. Toujours dans la théorie analytique, mais en utilisant ses propres outils, plutôt que la théorie cohomologique de Kashiwara, les définitions de seconde microlocalisation et de microlocalisation d'ordre supérieur dans le cas lagrangien sont dues à J. Sjöstrand [42]. Un peu plus tard, et pour traiter le problème d'interaction non linéaire de trois ondes progressives, J.-M. Bony introduit la seconde microlocalisation  $C^\infty$  sur une variété lagrangienne. Ce type de théorie a par ailleurs eu d'autres extensions: théorie isotrope [28], microlocalisations  $C^\infty$  d'ordre supérieur associées à des métriques de J.-M. Bony et N. Lerner [15], deuxième microlocalisation simultanée de J.-M. Delort [22]. Le champ d'application de la seconde microlocalisation ne se limite d'ailleurs pas à la compréhension des phénomènes non-linéaires: cet outil s'est révélé très performant pour l'étude des phénomènes de diffraction d'ondes linéaires.

On se limitera ici au cas de la deuxième microlocalisation analytique à croissance sur la lagrangienne  $T_\Lambda^*$ , où  $\Lambda$  est une sous-variété de  $\mathbb{R}^n$ , en théorie de J. Sjöstrand.

Soit  $x = (x', x'') \in \mathbb{R}^n$ ,  $x' \in \mathbb{R}^{n'}$ ,  $x'' \in \mathbb{R}^{n''}$ ,  $n' + n'' = n$  et  $\Lambda$  la sous-variété d'équation  $x'' = 0$ . On note  $(x, \xi) = (x', x'', \xi', \xi'')$  les points du fibré cotangent  $T^*\mathbb{R}^n$ . Le fibré conormal à  $\Lambda$ , a pour équation  $T_\Lambda^* = \{x'' = 0, \xi' = 0\}$ . On note  $(x', \xi'', x'^*, \xi''^*)$  les points du fibré cotangent  $T^*(T_\Lambda^*)$ . Le fibré cotangent à  $\Lambda$ ,  $T^*\Lambda$  s'identifie à un sous-espace de  $T^*(T_\Lambda^*)$  par l'injection

$$(17) \quad T^*\Lambda \ni (x', x'^*) \mapsto (x', 0; x'^*, 0) \in T^*(T_\Lambda^*).$$

Soit  $f$  une distribution à support compact dans  $\mathbb{R}^n$ . Pour  $w \in \mathbb{C}^n$ ,  $\lambda \in [1, \infty[$ ,  $\mu \in ]0, \mu_0]$ ,  $\mu_0 < 1$  on pose

$$(18) \quad T^2 f(w, \lambda, \mu) = \int_{\mathbb{R}^n} e^{-(\lambda/2)\bar{\mu}^2(w-y)^2} T^1 f(y, \lambda) dy, \quad \bar{\mu}^2 = \mu^2/1 - \mu^2$$

où  $T^1 f(y, \lambda)$  est défini pour  $y \in \mathbb{C}^n$ ,  $\lambda \in [1, \infty[$  par

$$(19) \quad T^1 f(y, \lambda) = \int e^{-(\lambda/2)y'^2 - \lambda/2(y' - x')^2 - \lambda/2(y'' - x'')^2} f(x) dx.$$

La transformation  $f \mapsto T^1 f$  est une transformation de Fourier-Bros-Iagolnitzer usuelle [42] qui envoie le complexifié  $(T_\Lambda^*)^{\mathbb{C}}$  sur la section nulle. Les transformées  $T^1 f$  et  $T^2 f$  vérifient des estimations uniformes d'espaces de Sjöstrand

$$(20) \quad \begin{aligned} |T^1 f(y, \lambda)| &\leq \text{cte } \lambda^A e^{\lambda/2(\text{Im } y)^2}; \\ |T^2 f(w, \lambda, \mu)| &\leq \text{cte } \lambda^A \left(\frac{2\pi}{\lambda \bar{\mu}^2}\right)^{n/2} e^{(\lambda/2)\mu^2(\text{Im } w)^2}. \end{aligned}$$

Pour  $\alpha = (x', \xi''; x'^*, \xi''^*)$  on pose

$$(21) \quad \chi(\alpha) = (w', w'') \quad w' = x' - ix'^*, \quad w'' = -\xi'' + i\xi''^*.$$

(22) **Définition 1.** Soit  $\alpha_0 \in T^*(T_A^*)$  et  $f \in \mathcal{E}'(\mathbb{R}^n)$ . On dit que  $\alpha_0$  n'appartient pas au deuxième micro-support à croissance de  $f$  le long de  $T_A^*$ , et on note  $\alpha_0 \notin WF_{T_A^*}^2(f)$  ssi il existe  $W$  voisinage de  $w_0 = \chi(\alpha_0)$ ,  $A, B, C, \mu_0 > 0$  tels que

$$(23) \quad \begin{aligned} \forall \mu \in ]0, \mu_0], \quad \forall w \in W, \quad \forall \lambda, \lambda \mu^2 \geq 1 \\ \Rightarrow |T^2 f(w, \lambda, \mu)| \leq A \lambda^B e^{\lambda(\mu^2/2)[(\text{Im } w)^2 - c]}. \end{aligned}$$

Alors  $WF_{T_A^*}^2(f)$  est un fermé homogène de  $T^*(T_A^*)$  et si  $P$  est un opérateur différentiel sur  $\mathbb{R}^n$ , on a  $WF_{T_A^*}^2(P f) \subset WF_{T_A^*}^2(f)$ .

(24) **Proposition 2** [27]. Soit  $f \in \mathcal{E}'(\mathbb{R}^n)$ , dont la transformée de Fourier vérifie

$$(25) \quad \exists M, \delta > 0 \text{ t.q. } \int (1 + |\xi''|)^\delta |\widehat{f}(\xi', \xi'')| d\xi'' \leq cte (1 + |\xi'|)^M.$$

La trace  $f|_A$  est bien définie et vérifie

$$(26) \quad SS(f|_A) \subset WF_{T_A^*}^2(f) \cap T^* A.$$

Soit à présent  $u$  solution de (1), telle que ses données de Cauchy  $u_0$  et  $u_1$  soient conormales classiques sur  $A = T_V^*(\mathbb{R}^d)$  où  $V$  est une hypersurface analytique réelle. Si  $D$  est un diagramme, on définit  $\Lambda_{[D]}$  et  $\Lambda_{\{D\}}$  comme les lagrangiennes complexes naturellement associées aux distributions  $[D]$  et  $\{D\}$ , à savoir

$$(27) \quad \begin{aligned} \Lambda_{\{D\}} = \{ (z_0, z_1, \dots, z_N, \zeta_0, \zeta_1, \dots, \zeta_N), \zeta_0 = 0, \zeta_i = 0 \text{ si } i \notin J \\ t_i = 0, \quad (x_i, \xi_i) \in T_V^* \mathbb{C}^d \text{ ou } \xi_i = 0 \text{ si } i \in J \}. \end{aligned}$$

$$(28) \quad \begin{aligned} \Lambda_{[D]} = \{ (z_0, z_1, \dots, z_N, \zeta_0, \zeta_1, \dots, \zeta_N), \text{ tel qu'il existe} \\ \Xi_1, \dots, \Xi_N \text{ avec } (z_i - z_{f(i)}, \Xi_i) \in \Lambda_\square, \quad \zeta_0 = \sum_{f(i)=0} \Xi_i, \\ \zeta_i = -\Xi_i + \sum_{f(j)=i} \Xi_j \text{ si } i \neq 0 \} \end{aligned}$$

où  $\Lambda_\square = T_{\Gamma^{\mathbb{C}}}^* \cup \{\zeta = 0\}$ ,  $\Gamma^{\mathbb{C}}$  étant le complexifié du cône d'onde ( $T_0^* \subset T_{\Gamma^{\mathbb{C}}}^*$ ).

Si  $A$  est la diagonale de  $M \times M$ ,  $M = \mathbb{R}^{(1+d)(N+1)}$ , on montre alors (en utilisant des résultats de finitude en géométrie sous-analytique réelle):

$$(29) \quad \text{Proposition 3} [27]. \quad WF_{T_A^*}^2([D] \otimes \{D\}) \subset C_{(T_A^*)^c}(\Lambda_{[D]} \otimes \Lambda_{\{D\}}).$$

Ici  $C_A(A)$  désigne le cône réel de  $A$  le long de  $A$ , sous ensemble du fibré normal à  $A$ , identifié au fibré conormal à  $A$  via la structure symplectique du fibré cotangent ambiant.

On déduit alors des Propositions 2 et 3 le résultat suivant, qui permet de rendre effectif le Théorème 5.

(30) **Proposition 4.**  $WF([D]\cdot\{D\}) \subset (\Lambda_{[D]}\widehat{\wedge} \Lambda_{\{D\}}) \cap T^*\mathbb{R}^{(1+d)(N+1)}$  où  $\widehat{\wedge}$  est l'opérateur de M. Kashiwara et P. Schapira [24]

$$(31) \quad (x, \xi) \in S_1 \widehat{\wedge} S_2 \text{ ssi il existe des suites } (x_n^i, \xi_n^i) \in S_i \\ \text{telles que } x_n^i \rightarrow x, \quad \xi_n^1 + \xi_n^2 \rightarrow \xi, \quad |x_n^1 - x_n^2| |\xi_n^i| \rightarrow 0.$$

#### IV. Estimations Géométriques des Singularités d'Ondes Non Linéaires

Soit  $u$  une onde semi-linéaire solution de (1) qui vérifie:

$$(32) \quad \begin{aligned} &\text{les données de Cauchy } u_0, u_1 \text{ sont conormales classiques} \\ &\text{sur } \Lambda = T_V^* \mathbb{R}^d \text{ où } V \text{ est une hypersurface analytique réelle.} \end{aligned}$$

Le Théorème 5 et la Proposition 4 fournissent des estimations géométriques des singularités de  $u$  à partir de l'ensemble des points limites d'ensembles de suites tracées dans le fibré cotangent complexe  $T^*\mathbb{C}^{1+d}$ , comme suit.

On note  $\mathcal{E}$  des ensembles de suites  $(z_n, \zeta_n) \in T^*\mathbb{C}^{1+d}$ ,  $z = (t, x)$ ,  $n \in \mathbb{N}$ , qui vérifient:

$$(33) \quad \text{la suite } z_n \text{ converge vers un point de } \Omega,$$

$$(34) \quad \begin{aligned} &\text{il existe une suite convergente } \eta_n \in \mathbb{C}^{1+d}, \quad |\eta_n| = 1 \text{ et} \\ &\text{une suite } \lambda_n \in \mathbb{C}^* \text{ telles que } \zeta_n = \lambda_n \eta_n, \end{aligned}$$

$$(35) \quad \begin{aligned} &\text{la suite } (z_n, \zeta_n) \text{ est caractéristique,} \\ &\text{i.e. } \zeta_n = (\tau_n, \xi_n), \quad \tau_n^2 = \xi_n^2. \end{aligned}$$

On suppose en outre que  $\mathcal{E}$  est stable par extraction de sous-suites, et vérifie:

$$(36) \quad \begin{aligned} &\mathcal{E} \text{ contient toute suite } (z_n, \zeta_n) \text{ satisfaisant} \\ &(33), (34), (35) \text{ et } \lim \zeta_n = 0, \end{aligned}$$

$$(37) \quad \begin{aligned} &\text{si } (z_n, \zeta_n) \in \mathcal{E} \text{ et } z'_n \text{ vérifie } \lim z'_n = \lim z_n \text{ et} \\ &\lim |z_n - z'_n| |\zeta_n| = 0 \text{ alors } (*) (z'_n, \zeta_n) \in \mathcal{E}, \end{aligned}$$

$$(38) \quad \begin{aligned} &\text{si } (z_n, \zeta_n) \in \mathcal{E} \text{ et si } z'_n \text{ est telle que } \lim z'_n \in \Omega \\ &\text{appartient au demi-cône d'onde de sommet } \lim z_n \\ &\text{qui ne rencontre pas l'hyperplan } t = 0 \text{ et si } (z_n, \zeta_n) \\ &\text{et } (z'_n, \zeta_n) \text{ sont sur la même bicaractéristique complexe} \\ &\text{de } \square, \text{ alors } (*) (z'_n, \zeta_n) \in \mathcal{E}. \end{aligned}$$

[(\*) signifie: il existe une suite extraite telle que.]

Soit à présent

$$(39) \quad \underline{\mathcal{E}} = \{\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_M, \dots\}$$

une suite (croissante) de tels ensembles de suites telles que

$$(40) \quad \text{si } (z_n, \zeta_n^j) \in \mathcal{E}_{k_j} \text{ sont } N \text{ suites } (j = 1, \dots, N)$$

possédant le même point de base  $z_n$  et si  $\zeta_n$  est une suite telle que  $(z_n, \zeta_n)$  vérifie (33), (34), (35) et  $\lim(\zeta_n^1 + \dots + \zeta_n^N - \zeta_n) = 0$  alors  $(*) (z_n, \zeta_n) \in \mathcal{E}_k$  avec  $k = k_1 + \dots + k_N$ .

On définit alors  $Z_M(\underline{\mathcal{E}})$  par :

$$(41) \quad \begin{aligned} Z_M(\underline{\mathcal{E}}) = & \{(z, \zeta) \in T^* \mathbb{C}^{1+d}|_{\Omega}, \text{ il existe } N \text{ suites} \\ & (z_n, \zeta_n^j) \in \mathcal{E}_{k_j}, z = \lim z_n, \zeta = \lim \zeta_n^1 + \dots + \zeta_n^N \\ & \text{et } k_1 + \dots + k_N \leq M\} \end{aligned}$$

[le point  $(z, \zeta)$  n'est pas caractéristique en général.]

Enfin on note  $\mathcal{A}_V$  l'ensemble des suites  $(z_n, \zeta_n)$  vérifiant (33), (34) et (35),  $z_n = (0, x_n)$ ,  $\zeta_n = (\tau_n, \xi_n)$ ,  $(x_n, \xi_n) \in T_{V^c}^*$  où  $V^c$  est un complexifié de  $V$ .

**Théorème 6.** Soit  $u$  une onde semi-linéaire solution de (1) avec  $s = (d/2) + \varrho$ ,  $\varrho > 0$ , vérifiant (32). Pour  $M \geq 1$  et  $\sigma \in [s + 1 + (M - 1)\varrho, s + 1 + M\varrho]$  on a  $WF^\sigma(u)|_{t>0} \subset Z_M(\underline{\mathcal{E}}) \cap T^*\Omega_+$  dès que  $\underline{\mathcal{E}} = (\mathcal{E}_1, \mathcal{E}_2, \dots)$  vérifie  $\mathcal{A}_V \subset \mathcal{E}_1$ .

On obtient alors le Théorème 4 comme conséquence du Théorème 6 dans [27] en construisant explicitement un  $\underline{\mathcal{E}}$  de la forme  $(\mathcal{E}_1, \mathcal{E}_1, \mathcal{E}_1, \dots)$  tel que  $Z(\underline{\mathcal{E}}) \cap T^*\Omega_+ = T_{s \cup Q_+}^*$  près de  $B$ .

Le Théorème 6 a été amélioré depuis par J.-M. Delort [22], [23].

## Références

- [1] S. Alinhac: Paracomposition et opérateurs paradifférentiels. Comm. in P.D.E. **11-4** (1986) 87–121
- [2] S. Alinhac: Evolution d'une onde simple pour des équations non-linéaires générales. Current Topics in P.D.E. Kinokuniya Tokyo 1986, pp. 63–90
- [3] S. Alinhac: Interaction d'ondes simples pour des équations complètement non-linéaires. Ann. Sci. Ec. Norm. Sup. 4ème série **21** (1988) 91–132
- [4] M. Beals: Self-spreading and strength of singularities for solutions to semi-linear wave equations. Ann. Math. **118** (1983) 187–214
- [5] M. Beals: Propagation of smoothness for nonlinear second order strictly hyperbolic equations. Proc. of Symposia in Pure Math. **43** (1985) 21–45
- [6] M. Beals et G. Métivier: Progressing wave solutions to certain nonlinear mixed problems. Duke Math. J. **53** (1986) 125–137
- [7] M. Beals et G. Métivier: Reflection of transversal progressing waves in nonlinear strictly hyperbolic mixed problems. Am. J. Math. **109** (1987) 335–360
- [8] J.-M. Bony: Calcul symbolique et propagation des singularités pour les équations aux dérivées partielles non linéaires. Ann. Sci. Ec. Norm. Sup., 4ème série **14** (1981) 209–246

- [9] J.-M. Bony: Propagation des singularités .... Sém. Goulaouic-Schwartz, Ec. Polytechnique, 1979–80, no. 22
- [10] J.-M. Bony: Interaction des singularités .... Sém. Goulaouic-Schwartz, Ec. Polytechnique, 1981–82, no. 2
- [11] J.-M. Bony: Interaction des singularités pour les équations de Klein-Gordon non linéaires. Sém. Goulaouic-Meyer-Schwartz, Ec. Polytechnique, 1983–84, no. 10
- [12] J.-M. Bony: Second microlocalization and interaction of singularities for non linear P.D.E.. Hyperbolic equations and related topics (Mizohata, ed.). Kinokuniya, 1986, pp. 11–49
- [13] J.-M. Bony: Singularités des solutions de problèmes hyperboliques non linéaires. Advances in Microlocal Analysis (Garnir, ed.), NATO ASI Series, vol. 168, Reidel, 1985, pp. 15–39
- [14] J.-M. Bony et N. Lerner: Quantification asymptotique .... Sém. E.D.P., Ec. Polytechnique, 1986–87, no. 2 et 3
- [15] J.-M. Bony et N. Lerner: Quantification asymptotique et microlocalisations d'ordre supérieur. Ann. Sci. Ec. Norm. Sup., 4ème série **22** (1989) 1–57
- [16] J.-Y. Chemin: Calcul paradifférentiel précisé et applications aux équations aux dérivées partielles non semilinéaires. Duke Math. J. **56**, no. 3 (1988) 431–469
- [17] J.-Y. Chemin: Interaction contrôlée dans les E.D.P. non linéaires strictement hyperboliques. Bull. Soc. Math. France **116** (1988) 341–383
- [18] J.-Y. Chemin: Interaction de trois ondes dans les équations semilinéaires strictement hyperboliques d'ordre 2. Comm. in P.D.E. **12** (11) (1987) 1203–1225
- [19] J.-Y. Chemin: Régularité de la solution d'un problème de Cauchy fortement non linéaire à données singulières en un point. Ann. Inst. Fourier **39** (1989)
- [20] J.-Y. Chemin: Evolution d'une singularité ponctuelle dans des équations strictement hyperboliques non linéaires. à paraître
- [21] F. David et M. Williams: Singularities of solutions to semilinear boundary value problems. Amer. J. Math. **109** (1987) 1087–1109
- [22] J.-M. Delort: Deuxième microlocalisation simultanée et front d'onde de produits. Ann. Sci. Ec. Norm. Sup., 4ème série **23** (1990) 257–310
- [23] J.-M. Delort: Conormalité des ondes semi-linéaires le long des caustiques. Séminaire E.D.P., Ec. Polytechnique, 1988–89, no. 15
- [24] M. Kashiwara, P. Shapira: Microlocal study of sheaves. Astérisque **128** (1985)
- [25] Y. Laurent: Théorie de la deuxième microlocalisation dans le domaine complexe. Progress in Mathematics, vol. 53. Birkhäuser, 1985
- [26] G. Lebeau: Problème de Cauchy semi-linéaire en 3 dimensions d'espace. Un résultat de finitude. J. Funct. Anal. **77** (1988)
- [27] G. Lebeau: Equations des ondes semi-linéaires II. Contrôle des singularités et caustiques non linéaires. Invent. math. (1989)
- [28] G. Lebeau: Deuxième microlocalisation sur les sous-variétés isotropes. Ann. Inst. Fourier Grenoble **35**, no. 2 (1985) 145–216
- [29] G. Lebeau: Front d'onde des fonctions non linéaires et polynômes. Séminaire E.D.P., Ec. Polytechnique, 1988–89, no. 10
- [30] E. Leichtnam: Régularité microlocale pour des problèmes de Dirichlet non linéaires non caractéristiques d'ordre deux à bord peu régulier. Bull. Soc. Math. France **115** (1987) 457–489
- [31] R. Melrose: Semi-linear waves with cusp singularities. Actes Journées E.D.P., St Jean de Monts 1987, no. 10
- [32] R. Melrose et N. Ritter: Interaction of progressing waves for semi-linear wave equation I. Ann. Math. **121** (1985) 149–236
- [33] R. Melrose et N. Ritter: Interaction of progressing waves for semi-linear wave equation II. Arkiv för Math. **25** (1987) 91–114

- [34] Y. Meyer: Remarques sur un théorème de J.-M. Bony. Suppl. di Rend. del Circolo Mat. di Palermo, 1981, pp. 1–20
- [35] Y. Meyer: Régularité des solutions des équations aux dérivées partielles non linéaires. Sémin. Bourbaki. Lecture Notes in Mathematics, vol. 842. Springer, Berlin Heidelberg New York 1980, pp. 293–302
- [36] A. Piriou: Calcul symbolique non linéaire pour une onde conormale simple. Ann. Inst. Fourier **38**, no. 4 (1988) 173–186
- [37] J. Rauch et M. Reed: Nonlinear microlocal analysis of semi-linear hyperbolic systems in one space dimension. Duke Math. J. **49** (1982) 397–475
- [38] J. Rauch et M. Reed: Singularities produced by the nonlinear interaction of three progressing waves; examples. Comm. in P.D.E. **7** (1982) 1117–1133
- [39] J. Rauch et M. Reed: Classical, conormal, semilinear waves. Séminaire E.D.P., Ec. Polytechnique, 1985–86, no. 5
- [40] N. Ritter: Progressing wave solutions to non-linear hyperbolic Cauchy problems. Ph. D. Thesis M.I.T., 1984
- [41] A. Sa Barreto: Interaction of conormal waves for fully semilinear wave equations. (à paraître)
- [42] J. Sjöstrand: Singularités analytiques microlocales. Astérisque **95** (1982)
- [43] M. Sablé-tougeron: Régularité microlocale pour des problèmes aux limites non linéaires. Ann. Inst. Fourier **36**, no. 1 (1986) 39–82
- [44] M. Williams: Spreading of singularities at the boundary in semilinear hyperbolic mixed problems I: microlocal  $H^{s,s'}$  regularity. Duke Math. J. **56** (1988) 17–40
- [45] M. Williams: Spreading of singularities at the boundary in semilinear hyperbolic mixed problems II: crossing and self-spreading. Trans. Amer. Math. Soc. **310** (1988)
- [46] C.J. Xu: Propagation au bord des singularités pour des problèmes de Dirichlet non linéaires d'ordre deux. Actes Journées E.D.P., St Jean de Monts, 1989, no. 20

# Static and Moving Defects in Liquid Crystals

*Fang Hua Lin*

Courant Institute of Mathematical Sciences, New York University  
New York, NY 10012, USA

## 1. Introduction

There have been many recent activities in the analysis of defects (or singularities) of solutions of partial differential equations. In these solutions defects often reveal crucial facets of certain nonlinear problems which they model. Here I shall survey some recent studies on static and moving defects in liquid crystals. One may find references [1], [2] and [3] useful for the discussion below.

Liquid crystals are optically anisotropic, even when they are at rest. Scientifically, defects of the optical director in liquid crystals have long been of interest. One can resolve, in experiment, individual defects and details of configurations near them, with relatively simple optical (polarizing) microscopes. Often, such configurations are not static but, in many cases, they change very slowly with time. Thus it is useful to first consider defects in static liquid crystals.

In the classical Oseen-Frank model, energy minimizing static configurations of liquid crystals can be described by unit vector fields on a 3-dimensional domain. They are closely related to the theory of harmonic maps into the sphere. The latter also provides us with precise information concerning those isolated point defects such as bounds on the number of point defects and behavior of configurations near each such isolated point.

It was observed, at least experimentally [4], that line and surface defects do occur in liquid crystals. Following the general order-parameter theory of Ginzburg-Landau, J. Ericksen posed a new mathematical model to tackle these phenomena. It turns out the study of Ericksen's model is related to the study of harmonic maps to singular spaces (in this case the singular spaces are circular cones in  $\mathbb{R}^4$  or Minkowski space  $\mathbb{R}^{3,1}$ ). Problems of harmonic maps into singular spaces arise also in the study of super rigidity and other geometrical or topological problems for which we refer to a recent work of M. Gromov and R. Schoen [5].

In studying line and surface defects, we introduced a new mathematical device which was based on H. Federer's dimension reduction principle (see [6] and [3]). It is a useful tool also to study level sets of solutions to elliptic and parabolic equations [7]. Defect sets can be characterized as preimages of the vertex of the circular cone under these maps, or equivalently, the vanishing sets of the orientational order of liquid crystals. We can use this device to estimate the Hausdorff dimension, as well as the Hausdorff measure of defects of static and moving liquid crystals. Moreover, local behavior of liquid crystal configurations near these defects may also be described rather precisely.

## 2. Oseen-Frank Model and Point Defects

**2.1** In the Oseen-Frank model, the energy minimizing static configurations of liquid crystals can be described by maps  $n : \Omega \subseteq \mathbb{R}^3 \rightarrow \mathbb{S}^2$  which minimize the energy functional:

$$\int_{\Omega} W(n, \nabla n) dx , \quad (2.1)$$

where

$$\begin{aligned} W(n, \nabla n) = & k_1 |\operatorname{div} n|^2 + k_2 (q + n \cdot \operatorname{curl} n)^2 + k_3 |n \wedge \operatorname{curl} n|^2 \\ & + (k_2 + k_4) [\operatorname{tr}(\nabla n)^2 - (\operatorname{div} n)^2] , \end{aligned} \quad (2.2)$$

the  $k_i$ 's and  $q$  are material constants. The defect set is defined to be the discontinuity set of the map  $n$ .

It was shown by Hardt, Kinderlehrer and myself [8] that for any bounded Lipschitz domain  $\Omega$  in  $\mathbb{R}^3$  and any  $n_0 \in H^{1/2}(\partial\Omega, \mathbb{S}^2)$ , there is a minimizer  $n$  of (2.1) with  $n = n_0$  on  $\partial\Omega$  provided that  $k_1, k_2, k_3 > 0$ . Moreover,  $n$  satisfies the following:

(i) for any compact  $K \subset \Omega$ ,

$$\int_K |\nabla n|^2 dx \leq c(K, \Omega, k_1, k_2, k_3) ; \quad (2.3)$$

(ii)  $\nabla n \in L_{loc}^q(\Omega)$  for some  $q > 2$ ;

(iii) the defect set  $\Sigma$  of  $n$  has Hausdorff dimension strictly less than one, and  $n$  is analytic in  $\Omega \setminus \Sigma$ . (See [9].)

In the special case  $k_1 = k_2 = k_3 = 1$  and  $k_4 = q = 0$ ,

$$W(n, \nabla n) = |\nabla n|^2 \quad (2.4)$$

which is the integrand for harmonic maps from  $\Omega$  into  $\mathbb{S}^2$ . Schoen and Uhlenbeck [10] have shown, in this case, that the defect set of  $n$  consists of isolated points. Moreover, when  $\partial\Omega$  and  $n_0 : \partial\Omega \rightarrow \mathbb{S}^2$  are smooth,  $n$  is also smooth near  $\partial\Omega$ . This is unknown for minimizers of (2.1).

**2.2** Much has been learned recently about defects of an energy minimizing harmonic map from a 3-dimensional domain to the sphere  $\mathbb{S}^2$ . First of all, at each such isolated defect, there is a unique tangent map, which follows from a general theorem of L. Simon [11]. Second, these energy minimizing tangent maps are classified by Brezis-Coron and Lieb [12]. They are of the form  $\pm R \circ \frac{x}{|x|}$  for some rotation  $R$  of  $\mathbb{R}^3$ . Moreover, by a theorem of L. Simon [13] and Gulliver-White [14], there are two positive constants  $C$  and  $\alpha$  (independent of maps) such that if  $n$  is energy minimizing from  $\mathbb{B}^3$  to  $\mathbb{S}^2$  and  $0$  is a defect of  $n$ , then

$$|n(x) - \phi\left(\frac{x}{|x|}\right)| \leq c|x|^\alpha \quad (2.5)$$

for some  $\phi(\cdot) = \pm R(\cdot)$ .

An interesting application of the classification of energy minimizing tangent maps is the following a priori estimate on the number of defect points which is shown by Almgren-Lieb [15] and, independently, by Hardt and myself [16]:

For each compact  $K \subset B^3$  there is a universal constant  $C(K)$  such that

$$\begin{aligned} \# \text{ of defects of } n \text{ in } K &\leq C(K) \\ \text{for any energy minimizing map } n \text{ from } \mathbb{B}^3 \text{ to } \mathbb{S}^2. \end{aligned} \quad (2.6)$$

Related to (2.6) is the following stability theorem [16]:

If  $g : \mathbb{S}^2 \rightarrow \mathbb{S}^2$  is such that  $\|g - id\|_{C^1(\mathbb{S}^2)} \leq \varepsilon_0$  (for some  $\varepsilon_0 > 0$ ), then any energy minimizing map  $n : B^3 \rightarrow \mathbb{S}^2$  with  $n = g$  on  $\mathbb{S}^2$  has a unique defect point  $a$  so that

$$\|n(x) - R_a \circ \frac{(x-a)}{|x-a|}\|_{C^\alpha(B^3)} \leq C\varepsilon_0^{1/4} \quad (2.7)$$

and

$$|a| \leq C\varepsilon_0^{1/2}, \quad \|R_a - id\| \leq C\varepsilon_0^{1/4} \quad (2.8)$$

for some rotation  $R_a$  of  $\mathbb{R}^3$  and for some positive constants  $C$  and  $\alpha$ .

The classification of energy minimizing tangent maps in higher dimensions or target spheres with nonstandard metrics remains as a difficult open problem.

When the domain is four-dimensional and the target is  $\mathbb{S}^2$ , one deduces from a recent theorem of L. Simon [17] and Hardt and myself [18] the following:

**Theorem A.** Let  $u : B^4 \rightarrow \mathbb{S}^2$  be an energy minimizing map, then the defect set of  $u$  is locally a finite union of a finite set and a finite family of  $C^{1,\alpha}$  curves with finitely many crossings. Moreover, the 1-dimensional Hausdorff measure of the defect set is locally finite.

**2.3** Continuous non energy minimizing harmonic maps are of interest: from both analysis and differential geometry aspects. For example a classical problem is to represent a homotopy class of maps between compact Riemannian manifolds by harmonic maps (see [19], [20]). Some partial results related to the theory of liquid crystals have been found by Bethuel, Brezis and Coron [21] and Giacinta, Monica and Souček [22].

It should be noted also that solutions to harmonic map systems with the Dirichlet boundary condition are not unique (see [21], [23] and [24]), and the defect sets of these maps can be much more complicated [24].

### 3. Line and Surface Defects

**3.1** To explain line and surface defects in liquid crystals, one uses Ericksen's model. In this new models, the energy minimizing configuration of liquid crystals is described by a pair,  $(s, n)$ , where  $s : \Omega \rightarrow [-1/2, 1]$  is a real function which denotes the variable degree of orientation and  $n : \Omega \rightarrow \mathbb{S}^2$  denotes the axis of optical director. Thus  $(s, n)$  minimizes a bulk-energy functional which, for particular choices of material constants involved, reduces to

$$\int_\Omega [k|\nabla s|^2 + s^2|\nabla n|^2 + \psi(s)] dx \quad (3.1)$$

with  $k > 0$ , where the potential function  $\psi$  is a positive  $C^2$ -function defined on  $(-1/2, 1)$  (cf. [25] and [3]).

Since  $s$  may be zero somewhere in  $\Omega$ , (3.1) is a rather degenerate variational integral. However, the change of variables  $u = sn$  reduces (3.1) to

$$\int_{\Omega} [(k-1)|\nabla s|^2 + |\nabla u|^2 + \psi(s)] dx . \quad (3.2)$$

The variational integral (3.2) is, essentially, the energy of the map  $(s, u) : \Omega \rightarrow \mathbb{C}_k$  where  $\mathbb{C}_k$  is the circular cone  $\{(t, u) \in \mathbb{R} \times \mathbb{R}^3 : |t| = \sqrt{k-1}|u|\}$ , for  $k \geq 1$ , in the Euclidian space  $\mathbb{R}^4$ , or, the circular cone  $\{(t, u) \in \mathbb{R} \times \mathbb{R}^3 : |t| = \sqrt{1-k}|u|\}$ , for  $k < 1$ , in the Minkowski space  $\mathbb{R}^{3,1}$ . Now the direct method in the calculus of variations implies the existence of a minimizer of (3.2) under the Dirichlet boundary condition that  $(s, t)|_{\partial\Omega} \in H^{1/2}(\partial\Omega, \mathbb{C}_k)$ . Moreover, when  $0 < k \leq 1$ , and  $\psi \equiv 0$ , the minimizer of (3.2) is unique (see [25]).

Regarding the regularity of the map  $(s, u)$ , one has the following result, see [25]:

**Theorem B.** *Let  $\Omega$  be a bounded  $C^{1,\alpha}$  domain in  $\mathbb{R}^3$  and let  $(s_0, u_0) \in C^{1,\alpha}(\partial\Omega, \mathbb{R}^4)$ , with  $|s_0| = |u_0|$  on  $\partial\Omega$ . Suppose  $(s, u)$  is a minimizer of (3.2) which satisfies the constraint  $|s| = |u|$  a.e. in  $\Omega$ . Then  $(s, u) \in C^\beta(\overline{\Omega})$ . Moreover, when  $0 < k < 1$ , both  $s$  and  $u$  are lipschitz continuous in  $\overline{\Omega}$ .*

It should be pointed out that  $(s, u)$  is analytic (if  $\psi$  is) or smooth (if  $\psi$  is) on the open set  $\{x \in \Omega : s(x) > 0\}$ . This follows from standard elliptic regularity theory.

We also note that the existence and partial regularity of minimizers of (3.2) were also established by L. Ambrosio in [26] and [27].

**3.2** Having seen the regularity of the minimizers of (3.2), one is then interested in the defect sets of the optical director  $n$ . It is shown in [3] and [25] that defect sets are precisely the nodal set of the orientational order  $s$ , i.e., the set  $\{x \in \Omega : s(x) = 0\}$ .

**Theorem C.** *Let  $(s, u)$  be a minimizer of (3.2). Then the set  $\{s = 0\}$  is either all of  $\Omega$  or is of Hausdorff dimension  $\leq 2$ . If, in addition,  $k > 1$  and  $s \not\equiv 0$ , then the set  $\{s = 0\}$  has the Hausdorff dimension  $\leq 1$ .*

One notices that, for any  $0 < k < 1$ , examples of minimizers  $(s, u)$  of (3.2) with  $s \geq 0$  and with  $\{s = 0\}$  being 2-dimensional were explicitly constructed in [28].

The proof of *Theorem C* is based on the dimension reduction principle of H. Federer [6] and the monotonicity of the function  $N(r)$  defined below (see also [3], [7] and [25]).

For  $a \in \Omega$ ,  $r \in (0, d_a)$ ,  $d_a = \text{dist}(a, \partial\Omega)$ ,  $B_r(a) = \{x \in \Omega : |x - a| < r\}$ , we define

$$\begin{aligned} D(r) &= \int_{B_r(a)} [(k-1)|\nabla s|^2 + |\nabla u|^2 + s \cdot \psi'(s)] dx , \\ H(r) &= \int_{\partial B_r(a)} (k-1)|s|^2 + |u|^2 \end{aligned} \quad (3.3)$$

and

$$N(r) = \frac{r D(r)}{H(r)} \quad (3.4)$$

provided that  $H(r) \neq 0$ . Then one has the following

**Lemma D.** *There are two positive constants  $r_0$  and  $C$  depending only on  $\psi$  such that the function*

$$N(r) e^{Cr} \quad (3.5)$$

*is a monotone increasing function of  $r \in (0, r_a)$ ,  $r_a = \min[d_a, r_0]$ .*

The monotonicity of  $N(r)$  is also a useful fact in the work [5]. In fact, suppose  $u$  is a harmonic map from  $\mathbb{R}^n$  to  $N$  with curvature of  $N \leq 0$ . Then the function

$$N(r) = \frac{r \int_{B_r(a)} |\nabla u|^2 dx}{\int_{\partial B_r(a)} d^2(0, u)} \quad (3.6)$$

is a monotone increasing function of  $r$ . Where  $d(0, u)$  is the intrinsic distance from a fixed point  $0$  to  $u$  in  $N$ . The proof of (3.6) is based on the monotonicity formula for energy and the fact that

$$\Delta d^2(0, u) \geq 2|\nabla u|^2 \quad (3.7)$$

for any harmonic maps  $u : \mathbb{R}^n \rightarrow N$ .

As a consequence of (3.6)  $u$  is locally uniformly lipschitz continuous (independent on  $N$ ).

Finally we also note that, by combining Lemma D and [18], one can estimate 2-dimensional (or 1-dimensional in some special cases) Hausdorff measure of  $\{s = 0\}$  and describe the structure of  $\{s = 0\}$ .

## 4. Moving Defects

**4.1** Equations adequate for the treatment of both static and dynamic phenomena in liquid crystals, called the Leslie-Ericksen theory, were developed during the 1960's [1]. In [29], Erickson derived a full set of dynamic equations for nematic liquid crystals with variable degree of orientation. Since motions of liquid crystals are generally slow, the motion of the optical directors is our main concern. After neglecting the small velocity of the fluid, the motion of the optical director can be described as evolution of harmonic maps from  $\Omega$  to  $\mathbb{C}_k$  (a circular cone in  $\mathbb{R}^4$  or  $\mathbb{R}^{3,1}$ ). More precisely, we let  $v = (s, u) : \Omega \rightarrow \mathbb{C}_k$  and let

$$\tau(v) = 0 \quad \text{in } \Omega \quad (4.1)$$

be the equations for harmonic maps from  $\Omega$  to  $\mathbb{C}_k$ . Then the evolution of the optical director satisfies

$$\frac{\partial}{\partial t} v = \tau(v), \quad \text{for } (x, t) \in \Omega \times (0, \infty) \quad (4.2)$$

with initial Cauchy data

$$v(x, 0) = v_0(x), \quad x \in \Omega \quad (4.3)$$

and Dirichlet boundary condition

$$v(x, t) = v_0(x) \quad \text{for } x \in \partial\Omega. \quad (4.4)$$

The problem (4.2), (4.3), (4.4) was first studied by Eells-Sampson [20] and more recently by Struwe and Chen [30] and many others. The main difference between our problem and those studied earlier is that the target manifold contains singularities. Nevertheless one has the following

**Theorem E.** *If  $v_0$  is uniformly lipschitz continuous on  $\overline{\Omega}$  and  $v_0|_{\partial\Omega} \in C^{1,\alpha}(\partial\Omega)$ , then (4.2), (4.3), (4.4) has a global weak solution  $v \in L^\infty((0, \infty), H^1(\Omega))$ . Moreover, for  $k \leq 1$ ,  $v$  is unique and satisfies*

$$\sup_{\Omega} |\nabla_x v|(\cdot, t) \leq c \sup_{\Omega} |\nabla_x v_0(x)|, \quad (4.5)$$

and for  $k > 1$ ,

$$\|v\|_{C^k(\overline{\Omega})}(\cdot, t) \leq c \sup_{\Omega} |\nabla_x v_0(x)|. \quad (4.6)$$

Here  $c$  and  $\beta$  are constants depending only on  $\Omega$  and  $k$ .

The proof of (4.5) is based on the observation that  $\mathbb{C}_k$  in  $\mathbb{R}^{3,1}$  is a negatively curved Riemannian submanifold. Even when  $k > 1$ ,  $s^2$  is a strictly convex function on  $\mathbb{C}_k$ . This geometrical fact can be used to show (4.6).

**4.2** Since the moving defects are precisely the set  $\{(x, t) : s(x, t) = 0\}$  in  $\Omega \times \mathbb{R}^+$ , the problem reduces to studying the nodal set of a solution to certain parabolic systems. In general it is still an interesting subject for future studies. However, there are several recent works by C. Fefferman and H. Donnelly [31] and by Hardt and Simon [32] concerning nodal sets of solutions to elliptic equations. In [7] we studied the corresponding question for a class of heat equations. Interestingly enough, the function  $N(r)$  introduced in (3.4) can also be used to control not only the local behavior of solutions near nodal sets but also global Hausdorff measure of nodal sets. Generalizing those arguments in [7] to the problem (4.2) will be an interesting problem for future researchers.

## References

1. Ericksen, J. L., Kinderlehrer, D.: Theory and applications of liquid crystals. IMA vol. 5. Springer, Berlin Heidelberg New York 1986
2. Brezis, H.:  $\mathbb{S}^k$ -valued maps with singularities. (Lecture Notes in Mathematics, vol. 1365.) Springer, Berlin Heidelberg New York 1989
3. Lin, F.-H.: Nonlinear theory of defects in nematic liquid crystals; phase transition and flow phenomena. CPAM, vol. XLII, 1989, pp. 789–8914
4. Kléman, M.: Points, lines and walls. Wiley, New York 1983
5. Gromov, M., Schoen, R.: Harmonic maps into singular spaces and superrigidity. Preprint
6. Federer, H.: The singular sets of area-minimizing rectifiable currents with codimension one and of area-minimizing flat chain modulo two with arbitrary codimension. Bull. Amer. Math. Soc. 76 (1970) 667–711

7. Lin, F.-H.: Nodal sets of solutions of elliptic and parabolic equations. To appear in CPAM
8. Hardt, R., Kinderlehrer, D., Lin, F. H.: Existence and partial regularity of static liquid crystal configurations. *Comm. Math. Physics* **105** (1986) 547–570
9. Hardt, R., Kinderlehrer, D., Lin, F. H.: Stable defects of minimizers of constrained variational principles. *Ann. Inst. H. Poincaré, Anal. Nonlinéaire* **5**, no. 4 (1988) 297–322
10. Schoen, R., Uhlenbeck, K.: A regularity theory of harmonic maps. *J. Diff. Geom.* **18** (1983) 253–268
11. Simon, L.: Asymptotics for a class of nonlinear evolution equations. *Anal. Math.* **118** (2) (1983) 525–571
12. Brezis, H., Coron, J. M., Lieb, E.: Harmonic maps with defects. *Comm. Math. Phys.* **107** (1986) 649–705
13. Simon, L.: Isolated singularities for extreme of geometrical variational problems. (*Lecture Notes in Mathematics*, vol. 1161.) Springer, Berlin Heidelberg New York 1986
14. Gulliver, R., White, B.: The rate of convergence of a harmonic map at a singular point. *Math. Ann.* **283** (1989) 539–550
15. Almgren, F. Jr., Lieb, E.: Singularities of energy minimizing maps from the ball to the sphere: examples, counterexamples, and bounds. *Anal. Math.* **128** (1988) 483–530
16. Hardt, R., Lin, F. H.: Stability of singularities of minimizing harmonic maps. *J. Diff. Geom.* **29** (1988) 113–123
17. Simon, L.: On the singularities of harmonic maps. Preprint
18. Hardt, R., Lin, F. H.: The singular set of an energy minimizing map from  $\mathbb{B}^4$  to  $\mathbb{S}^2$ . Preprint
19. Eells, J., Lemaire, L.: Another report on harmonic maps. *Bull. London Soc.* **20** (8) (1988) Part 5
20. Eells, J., Sampson, J. H.: Harmonic mappings of Riemannian manifolds. *Amer. J. Math.* **86** (1964) 109–160
21. Bethuel, F., Brezis, H., Coron, J. M.: Relaxed energies for harmonic maps. Preprint
22. Giaquinta, M., Monica, G., Souček, J.: Cartesian currents and liquid crystals, dipoles, singular lines and singular points. Preprint
23. Hardt, R., Kinderlehrer, D., Lin, F. H.: The variety of configurations of static liquid crystals. To appear
24. Hardt, R., Poon, C. C., Lin, F. H.: Axially symmetric harmonic maps minimizing a relaxed energy. Preprint
25. Lin, F. H.: On nematic liquid crystals with variable degree of orientations. To appear in CPAM
26. Ambrosio, L.: Existence of minimal energy configurations of nematic liquid crystals with variable degree of orientation. Preprint
27. Ambrosio, L.: Regularity of solutions of a degenerate elliptic variational problem. Preprint
28. Ambrosio, L., Virga, E.: A boundary value problem for nematic liquid crystals with variable degree of orientations. Preprint
29. Ericksen, J. L.: Liquid crystals with variable degree of orientation. IMA preprint # 559 (1989)
30. Struwe, M.: The evolution of harmonic maps. In these proceedings, pp. 1197–1203
31. Donnelly, H., Fefferman, C.: Nodal sets of eigenfunctions on Riemannian manifolds. *Inv. Math.* **93** (1988) 161–183
32. Hardt, R., Simon, L.: Nodal sets for solutions of elliptic equations. *J. Diff. Geom.* **30** (1989) 505–522



# On Kinetic Equations \*

Pierre-Louis Lions

CEREMADE, Université Paris-Dauphine, Place de Lattre de Tassigny  
F-75775 Paris Cedex 16, France

## I. Introduction

We will review some recent progress on various *kinetic equations* which include the well-known *Boltzmann equation* and also *Vlasov models* like Vlasov-Poisson or Vlasov Maxwell systems. Before describing more precisely these mathematical results, we would like to recall first a few of the basic physical notions underlying these models.

First of all, kinetic models are a branch of *Statistical Physics*. And they arise in a large number of different physical contexts like, for instance, in the study of the dynamics of electrons or ions in plasmas, in the study of the dynamics of nucleons in Nuclear Physics, in the modelling of semi-conductors, in the modelling of the reentry of various aircrafts in a rarefied atmosphere, in the study of the formation and stability of planetary rings or even in the study of the formation of galaxies. It is worth noting that this list, by no means exhaustive, includes different physical interactions on very different scales. Of course, each of these applications reveals specific mathematical questions that we will not address here. Instead, we will concentrate here on the main mathematical issues raised by all these applications.

In spite of these extremely different physical backgrounds, the main principle underlying these models can be summarized as follows. Let us suppose we want to study the evolution of a large number of particles that we take to be identical in order to simplify the presentation. The reader should be aware that the word particle above is used here in a vague sense and might be rather misleading since in the applications listed above the “particles” may be electrons, ions, nucleons, molecules, rocks or stars ! Therefore, in order to be more precise, we consider these objects as classical point-particles.

Next, we observe that when we deal with a large number of particles, it is impossible to study the evolution of each particle and we wish to look instead for a statistical description of this evolution. In other words, the unknowns which would have been otherwise the positions and velocities of each particle “reduce” to a function  $f$  of  $(x, v, t)$ ,  $(x, v \in \mathbb{R}^N, t \geq 0)$  which is the density of particles

---

\* Dedicated to the memory of R.J. DiPerna.

at position  $x$ , time  $t$  and with velocity  $v$ . Of course,  $f$  is nonnegative. Next, we simply indicate that many kinetic models, but not all of them, take the following form

$$\frac{\partial f}{\partial t} + v \cdot \nabla_x f + F \cdot \nabla_v f = C \quad \text{in } \mathbb{R}^N \times \mathbb{R}^N \times (0, \infty). \quad (1)$$

Here and everywhere below,  $\nabla_x f$  and  $\nabla_v f$  denote respectively the gradient of  $f$  with respect to  $x$  and  $v$ . Variants involve different velocities than  $v$  or convolution terms instead of  $F \cdot \nabla_v f$  ... In (1),  $F$  stands for a *force* acting on the particles while  $C$  is a term which takes into account the possible *collisions* between particles.

Let us now give a few specific examples: we begin with collision-less models that is we set  $C = 0$ . These models are called *Vlasov models*. In the case of charged particles thus obeying the electromagnetic interaction, one uses two models namely the *Vlasov-Maxwell system* and the *Vlasov-Poisson system*. The Vlasov-Maxwell system consists of (1) (with  $C = 0$  and  $N = 3$ ) and of

$$F(x, v, t) = E(x, t) + v \times B(x, t) \quad \text{in } \mathbb{R}_x^3 \times \mathbb{R}_v^3 \times (0, \infty) \quad (2)$$

$$\begin{aligned} \frac{\partial E}{\partial t} - \frac{1}{c} \operatorname{curl} B &= -j, & \frac{\partial B}{\partial t} + \frac{1}{c} \operatorname{curl} E &= 0, \\ \operatorname{div} B &= 0, & \operatorname{div} E &= \varrho \quad \text{in } \mathbb{R}_x^3 \times (0, \infty) \end{aligned} \quad (3)$$

$$j_k = \int_{\mathbb{R}^3} f v_k dv, \quad \varrho = \int_{\mathbb{R}^3} f dv. \quad (4)$$

Observe that (2) means that  $F$  is the *Lorentz force* associated to the *electromagnetic field*  $(E, B)$  which of course satisfies the *Maxwell equations* (3). Finally, (4) means that the electromagnetic field is created by the particles. And this field creates a force ( $F$ ) which acts on the particles – this is why one speaks of self-consistent forces for Vlasov models.

The parameter  $c$  is the speed of light and one sees that, in the limit  $c \rightarrow \infty$ , (2)–(4) becomes

$$F(x, t) = E(x, t) = -\nabla_x \Phi(x, t) \quad \text{in } \mathbb{R}_x^3 \times (0, \infty) \quad (5)$$

$$-\Delta \Phi = \varrho \quad \text{in } \mathbb{R}_x^3 \times (0, \infty) \quad (6)$$

$$\varrho(x, t) = \int_{\mathbb{R}^3} f dv. \quad (7)$$

This system coupled with (1) (with  $C = 0$ ,  $N = 3$ ) is called the *Vlasov-Poisson system*.

It is important for the analysis we will present to observe two facts. First, in Vlasov models or more specifically in collisionless models ( $C = 0$ ), the equation (1) simply means that  $f$  is *constant along particle paths*. Of course, these paths are given by Newton's law ( $\dot{x} = v$ ,  $\dot{v} = F$ ). Next, a general feature of Vlasov models

illustrated by the two examples above is the particular dependence of the force upon  $f$ : indeed,  $F$  depends in fact only on *macroscopic quantities* that is averages of  $f$  in  $v$  like

$$\int_{\mathbb{R}^3} f(x, v, t) \psi(v) dv$$

for some given  $\psi$ .

We now give one example of *collision models* or in other words one example for the term  $C$ . This is a famous nonlinear term introduced by J.C. Maxwell [35] and L. Boltzmann [6] which is given by

$$C = Q(f, f) = \int_{\mathbb{R}^3} dv_* \int_{S^2} d\omega B(v - v_*, \omega) \{f' f'_* - f f_*\} \quad (8)$$

where  $f_* = f(x, v_*, t)$ ,  $f' = f(x, v', t)$ ,  $f'_* = f(x, v'_*, t)$  and we have

$$v' = v - (v - v_*, \omega)\omega, \quad v'_* = v_* + (v - v_*, \omega)\omega. \quad (9)$$

The collision operator  $C$  determines the rate at which the density  $f$  is modified by collisions taking place at  $(x, t)$ , between two particles with velocities  $v$  and  $v_*$  – in other words, one makes a statistical balance of all possible collisions affecting the density  $f(x, v, t)$ . Now, if we assume these collisions to be elastic, the velocities  $v', v'_*$  after the collision satisfy

$$v' + v'_* = v + v_*, \quad |v'|^2 + |v'_*|^2 = |v|^2 + |v_*|^2. \quad (10)$$

All possible solutions of (10) are given by (9), where  $\omega$  is a parameter allowing a simple description of the set of solutions of (10).

Finally,  $B$  is always assumed to be a nonnegative function of  $|v - v_*|$  and  $(v - v_*, \omega)$  only. It depends on the interactions between the particles and the simplest (and most famous) example is the so-called *hard-spheres* case where  $B(z, \omega) = |(z, \omega)|$ .

Let us mention that the *Boltzmann equation* is the equation (1) with  $F = 0$  and  $C$  given as above.

Other models for the collision term  $C$  exist: some are modifications of the Boltzmann model like the so-called Boltzmann-Enskog or Boltzmann-Dirac models while others are, in some vague sense, simplifications like the Fokker-Planck or the Landau models. We refer to C. Cercignani [8], S. Chapman and T.G. Cowling [10], C. Truesdell and R.G. Muncaster [40] for more details and also for derivations of the Boltzmann model from first principles.

After this brief description of three famous examples of kinetic models, we now turn to a quick presentation of the numerous mathematical problems raised by these models. Let us mention immediately that we shall concentrate on the main mathematical issues forgetting many specific questions of interest for one or several of the physical applications listed above.

The first category of problems is the study of the Cauchy problems for these models, prescribing  $(f, E, B)$  (or  $f$  depending on the model) at time  $t = 0$  with the usual questions regarding existence, uniqueness, regularity, approximation and

numerical analysis, special solutions, steady states, stability, long-time behavior, boundary conditions ...

But there is much more at hand. In fact, it seems fair to say that Boltzmann's equation is not only famous because of its fascinating mathematical structure but also because of its formal relations with other famous mathematical physics models like hydrodynamical models (compressible Euler equations i.e. gaz dynamics systems, incompressible Navier-Stokes or Stokes or Euler equations ...). It is worth recalling Hilbert's goals to solve Boltzmann equation and to recover from it the Fluid Dynamics equations. Therefore, an important category of mathematical problems concerns the systematic study of the numerous links between kinetic models and other models in Physics. We just mentioned hydrodynamical limits and the link with Fluid Dynamics but one has to add to that theme the derivation of MHD equations, combustion models, reaction-diffusion systems and in fact hyperbolic systems of conservation laws. Indeed, arbitrary (symmetric) systems of conservation laws can be formally approximated by ad hoc kinetic models.

Other limits involve the derivation of kinetic models from "large number of particles limits" and the study of statistical solutions or hierarchies of equations and the propagation of chaos. Another example of a connection with another physical regime is the derivation of kinetic models from Quantum models via Wigner transforms and semi-classical limits ...

We will concentrate here on the first category of problems namely the analysis of the Cauchy problem for kinetic models even if it is quite clear that progress on this theme should and has already yielded progress on the second theme as well. And even if we restrict our attention to the basic existence and uniqueness questions, three types of results have been obtained:

- smooth and unique solutions in the small that is locally in time or globally with smallness restrictions: for the Boltzmann equation alone, many important works have been given in that direction and we can quote only a few of them (complete lists of references can be found in the references given in the bibliography here) like H. Grad [27], C. Cercignani [8], R. Illner and M. Shinbrot [30], T. Nishida and K. Imai [36], S. Ukaï [41], K. Hamdache [28]
- ...
- existence and uniqueness of special solutions: a famous example is given by the study of space-homogeneous solutions, i.e.  $x$ -independent solutions of Boltzmann equation, study initiated by the work of T. Carleman [7] ...
- global existence of weak solutions.

This last category of results has been obtained in a series of works for all kinetic models by R.J. DiPerna and the author [14–18] (see also the survey [19]). And we are going to give one sample of these results by specializing our attention to the Vlasov-Poisson-Boltzmann model in the next section. Finally, in the last section, we will indicate what are the main tools of independent interest that are being used in the proofs of these results.

## II. Global Existence of Weak Solutions

As we said before, we concentrate on the Vlasov-Poisson-Boltzmann (VPB in short) system:

$$\frac{\partial f}{\partial t} + v \cdot \nabla_x f - \lambda \nabla_x \Phi \cdot \nabla_v f = Q(f, f) \quad \text{in } \mathbb{R}_x^3 \times \mathbb{R}_v^3 \times (0, \infty) \quad (11)$$

$$-\Delta \Phi = \varrho \quad \text{in } \mathbb{R}_x^3 \times (0, \infty) \quad (12)$$

and in order to make  $\Phi$  well-defined by (12) we prescribe  $\Phi = 0$  at infinity in a weak sense (like for instance  $\Phi \in L^\infty(0, \infty; L^6(\mathbb{R}_x^3))$ ). And we recall

$$\begin{aligned} Q(f, f) &= Q^+(f, f) - Q^-(f, f), \\ Q^-(f, f) &= f L(f), \quad L(f) = f *_{_v} A \end{aligned} \quad (13)$$

$$\begin{aligned} Q^+(f, f) &= \int_{\mathbb{R}^3} dv_* \int_{S^2} d\omega B(v - v_*, \omega) f' f'_*, \\ Q^-(f, f) &= \int_{\mathbb{R}^3} dv_* \int_{S^2} d\omega B(v - v_*, \omega) f f_*, \end{aligned} \quad (14)$$

$$v' = v - (v - v_*, \omega)\omega, \quad v'_* = v + (v - v_*, \omega)\omega. \quad (16)$$

Finally,  $\lambda$  is a nonnegative parameter: when  $\lambda = 0$ , the (VPB) system reduces to the standard Boltzmann equation, while when  $\lambda > 0$  and  $B \equiv 0$  the (VPB) system becomes the Vlasov-Poisson system.

We will make the following assumptions

$$B(z, \omega) \text{ is a function of } z \text{ and } (z, \omega) \text{ only, } B \geq 0 \quad (17)$$

$$\int_{|z| \leq R} \int_{S^2} B dz d\omega < \infty, \quad \frac{1}{1 + |z|^2} \int_{|v-z| \leq R} \int_{S^2} B dv d\omega \rightarrow 0, \quad (18)$$

as  $|z| \rightarrow \infty$ , (for all  $R < \infty$ ).

These assumptions are clearly satisfied in the case of the hard-spheres model  $B(v, \omega) = |(v, \omega)|$ . Let us also mention that (18) corresponds to the classical angular cut-off assumption.

Let us now recall the known a priori estimates (when  $B \not\equiv 0$ ) which have all a physical origin (conserved quantities or decay of entropy): formally, a solution of (11)–(12) satisfies

$$\begin{aligned} \iint_{\mathbb{R}^3 \times \mathbb{R}^3} f \psi dv dx \quad &\text{is independent of } t \geq 0, \\ \text{where } \psi = 1, v_\alpha \quad (\alpha = 1, 2, 3), \end{aligned} \quad (19)$$

$$\iint_{\mathbb{R}^3 \times \mathbb{R}^3} f|v|^2 dx dv + \lambda \int_{\mathbb{R}^3} |\nabla_x \Phi|^2 dx \text{ is independent of } t \geq 0, \quad (20)$$

$$\frac{d}{dt} \iint_{\mathbb{R}^3 \times \mathbb{R}^3} f|x - vt|^2 = -\lambda t \int_{\mathbb{R}^3} |\nabla_x \Phi|^2 dx. \quad (21)$$

Finally, the decay of (the mathematical) entropy which is a crucial part of the famous  $H$ -theorem is expressed by the following formal identity:

$$\begin{aligned} \frac{d}{dt} \iint_{\mathbb{R}^3 \times \mathbb{R}^3} f \log f dx dv + \frac{1}{4} \int_{\mathbb{R}^3} dx \iint_{\mathbb{R}^3 \times \mathbb{R}^3} dv dv_* \\ \cdot \int_{S^2} B(f'f'_* - ff_*) \log \frac{f'f'_*}{ff_*} d\omega = 0. \end{aligned} \quad (22)$$

Notice that the second term is clearly nonnegative.

From this, one deduces easily that, if the initial condition  $f^0$  satisfies for some  $R \in (0, \infty)$

$$\iint_{\mathbb{R}^3 \times \mathbb{R}^3} f^0(1 + |x|^2 + |v|^2 + |\log f^0|) dx dv + \lambda \int_{\mathbb{R}^3} |\nabla_x \Phi^0|^2 dx \leq R, \quad (23)$$

then, formally at least, a solution of (11)–(12) satisfies for all  $t \geq 0$

$$\left\{ \begin{array}{l} \iint_{\mathbb{R}^3 \times \mathbb{R}^3} f(1 + |v|^2) dx dv + \lambda \int_{\mathbb{R}^3} |\nabla_x \Phi|^2 dx \leq R \\ \iint_{\mathbb{R}^3 \times \mathbb{R}^3} f(|x|^2 + \log |f|) \leq K(1+t)^2 \\ \int_0^t ds \int_{\mathbb{R}^3} dx \iint_{\mathbb{R}^3 \times \mathbb{R}^3} dv dv_* \\ \int_{S^2} B(f'f'_* - ff_*) \log \frac{f'f'_*}{ff_*} d\omega \leq K(1+t)^2, \end{array} \right. \quad (24)$$

for some positive constant  $K$  which depends only on  $R$ . In fact, if  $\lambda = 0$  (the case of the pure Boltzmann equation), the factor  $(1+t)^2$  can be omitted in the last inequality. In the collisionless case (Vlasov-Poisson system) i.e. when  $B \equiv 0$ , the Liouville conservation is translated by the following a priori estimate (formally)

$$\iint_{\mathbb{R}^3 \times \mathbb{R}^3} \beta(f) dx dv \text{ is independent of } t \quad (25)$$

whenever  $\beta \in C([0, \infty); [0, \infty))$  and  $\beta(f^0) \in L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3)$ .

However, all these bounds are not enough to allow to define  $Q(f, f)$  in a meaningful way. It only allows to define  $Q^\pm(f, f)$  as an a.e. finite function on  $(0, \infty) \times \mathbb{R}_x^3 \times \mathbb{R}_v^3$  ( $L^{1/2}$  in fact). In particular, it is not clear how one should define a solution of (11)–(12).

To circumvent this difficulty, we introduced in [14, 15] the notion of renormalized solutions. The first step is to prove that an additionnal bound can be obtained:

$$\int_0^M dt \int_{\mathbb{R}^3} dx \int_{|v| \leq M} \{Q^+(f, f) + Q^-(f, f)\} \frac{1}{1+f} dv \leq C \quad \text{for all } M > 0 \quad (26)$$

where  $C$  is a positive constant depending only on  $M$  and on  $R$  (in (23)).

In particular,  $Q^\pm(f, f)(1+f)^{-1}$  makes sense in  $L^1_{loc}$ . We may now give the

**Definition.**  $f \in C([0, \infty); L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3))$  is a renormalized solution of (11)–(12) if it satisfies (24), (26), (12) in the sense of distributions, and

$$\frac{\partial}{\partial t} \beta(f) - \operatorname{div}_x(v\beta(f)) - \lambda \operatorname{div}_v(\Phi\beta(f)) = Q(f, f)\beta'(f) \quad \text{in } \mathcal{D}' \quad (27)$$

for all  $\beta \in C([0, \infty); \mathbb{R})$  such that  $\beta'(t)(1+t)$  is bounded on  $[0, \infty)$ .

**Remarks.** 1) (27) is nothing else than writing formally the chain rule using (11).

2) It is enough to consider a single function  $\beta$  provided it is one to one, like for instance  $\beta(t) = \log(1+t)$ .

3) One can show (see [15]) that the notion of renormalized solution for Boltzmann equations ( $\lambda = 0$ ) is stronger than the more standard notion of mild solution (an integrated form of the equation along almost all particle paths).

4) It is worth recalling that in the collision-less case ( $B = 0$ ), the physical law behind the equation is that  $f$  should be constant along particle paths, a fact which is clearly equivalent to the fact that  $\beta(f)$  is constant along particle paths (for all  $\beta$  or for a single 1-1  $\beta \dots$ ). This is precisely what we request in the above definition.

5) In the collision-less case ( $B = 0$ ) i.e. in the case of the Vlasov-Poisson system, renormalized solutions are stronger than the classical weak solutions as built by Arsenev [4]. In particular (see [18]) they satisfy the Liouville conservation (25) and are continuous in time (with values in  $L^p_{x,v} \dots$ ).

Our main existence result is the

**Theorem 1.** Let (17)–(18) hold and let  $f^0$  satisfy (23), then there exists a renormalized solution  $f$  of (11)–(12) in  $C([0, \infty); L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3))$  which satisfies  $f|_{t=0} = f^0$  on  $\mathbb{R}_x^3 \times \mathbb{R}_v^3$  and

$$\begin{aligned} & \frac{d}{dt} \iint_{\mathbb{R}_x^3 \times \mathbb{R}_v^3} f \log f dx dv + \frac{1}{4} \int_{\mathbb{R}^3} dx \iint_{\mathbb{R}^3 \times \mathbb{R}^3} dv dv_* \\ & \cdot \int_{S^2} B(f' f'_* - f f_*) \log \frac{f' f'_*}{f f_*} d\omega \\ & \leq \iint_{\mathbb{R}_x^3 \times \mathbb{R}_v^3} f^0 \log f^0 dx dv. \end{aligned} \quad (28)$$

**Remarks.** 1) Uniqueness and regularity of solutions are major open problems. Also, the equality in (28) is another important open question. In some vague

sense, this result is analogous to the existing results (essentially due to J. Leray [31]) on three-dimensional incompressible Navier-Stokes equations.

2) Theorem 1 and its proof has been applied, adapted or extended to various other kinetic equations and related questions: K. Hamdache [29] (boundary conditions); B. Perthame [37] (BGK model); J. Polewczak [39], L. Arkeryd and C. Cercignani [3] (Boltzmann-Enskog model); M.J. Esteban and B. Perthame [23] (inelastic collisions); J.M. Dolbeault [13] (Boltzmann-Dirac model); L. Arkeryd [2], L. Desvillettes [11], C. Cercignani [9] (long time behavior); L. Desvillettes [12] (convergence of splitting methods); C. Bardos, F. Golse and D. Levermore [5] (convergence to some incompressible Fluid Dynamics models when the mean free path goes to 0).

The main ingredient in the proof of Theorem 1 is the following “stability under weak convergence” result taken from [15, 16].

**Theorem 2.** *Let  $f_n^0$  satisfy (23) and let (17)–(18) hold. Let  $f_n \in C([0, \infty); L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3))$  be a renormalized solution of (11)–(12) satisfying  $f_n|_{t=0} = f_n^0$ . Without loss of generality, we may assume that  $f_n^0, f_n$  converge weakly in  $L^1$  to  $f^0, f$  respectively. Then, we have*

- 1)  $\int_{\mathbb{R}_v^3} f_n \psi(v) dv \xrightarrow{n} \int_{\mathbb{R}_v^3} f \psi(v) dv$  in  $L^1((0, T) \times \mathbb{R}_x^3)$  for all  $T \in (0, \infty)$  and for all  $\psi \in L^\infty(\mathbb{R}_v^3)$ .
- 2)  $\int_{\mathbb{R}_v^3} Q^\pm(f_n, f_n) \psi(v) dv \xrightarrow{n} \int_{\mathbb{R}_v^3} Q^\pm(f, f) \psi(v) dv$  locally in measure on  $(0, \infty) \times \mathbb{R}_x^3$  for all  $\psi \in L^\infty(\mathbb{R}_v^3)$  with compact support.
- 3)  $\liminf_{t \rightarrow \infty} \iint_{\mathbb{R}_x^3 \times \mathbb{R}^3} dv dv_* \int_{S^2} B d\omega(f'_n f'_{n*} - f_n f_{n*}) \log \frac{f'_n f'_{n*}}{f_n f_{n*}} \geq \iint_{\mathbb{R}_x^3 \times \mathbb{R}^3} dv dv_* \int_{S^2} B d\omega(f' f'_* - f f_*) \log \frac{f' f'_*}{f f_*}$ , a.e.  $t, x$ .
- 4)  $f \in C([0, \infty); L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3))$  is a renormalized solution of (11)–(12) which satisfies  $f|_{t=0} = f^0$  on  $\mathbb{R}_x^3 \times \mathbb{R}_v^3$ .

As we said in the above remarks, uniqueness and regularity of solutions are not known in general. However, in the special case of the Vlasov-Poisson system that is when  $B \equiv 0$ , these questions are solved in [33] and follow from the following new a priori estimates.

**Theorem 3.** *Let  $B \equiv 0$  (Vlasov-Poisson system) and let  $f^0 \in L^1 \cap L^\infty(\mathbb{R}_x^2 \times \mathbb{R}_v^3)$  satisfy for some  $k > 3$*

$$\iint_{\mathbb{R}_x^3 \times \mathbb{R}_v^3} f^0 |v|^m dv dv < \infty \quad \text{for all } m \in [0, k]. \quad (29)$$

*Then, there exists a renormalized solution  $f \in C([0, \infty); L^1(\mathbb{R}_x^3 \times \mathbb{R}_v^3))$  satisfying*

$$\begin{aligned} \sup_{t \in [0, T]} \iint_{\mathbb{R}_x^3 \times \mathbb{R}_v^3} f(x, v, t) |v|^m dv dx &< \infty, \\ \text{for all } m \in [0, k], T \in (0, \infty). \end{aligned} \quad (30)$$

**Remarks.** 1) We prove in fact in [33] a bound on  $\iint_{\mathbb{R}^3 \times \mathbb{R}^3} f(x, v, t) \cdot |v|^m dv dx$  which depends only on the bounds on  $f^0$ .

2) Some uniqueness and regularity results have been independently proven by K. Pfaffelmoser [38].

### III. Mathematical Tools for the Global Existence of Weak Solutions in Kinetic Models

In this last section, we briefly present the three tools of independent interest that we use in the proofs of the global existence results (like Theorem 1) for weak solutions of kinetic models. These tools are 1) the notion of renormalized solutions and their properties, 2) a theory of a.e. flow-solutions of ordinary differential equations, 3) velocity-averaging lemmata.

We already described above the first tool (and its physical interpretation): we just want to mention here that this notion can be useful even for nonlinear equations involving second-order terms like Fokker-Planck-Boltzmann equations [14] or Landau equations [19] or like nonlinear elliptic equations (see P.L. Lions and F. Murat [32]). Let us emphasize the elementary fact that this notion consists essentially in writing down the equation formally satisfied by a nonlinear change of unknown.

The second tool is a theory of a.e. flows developed by R.J. DiPerna and the author [20, 21] that can be illustrated by the following example. Let  $\Omega$  be a bounded smooth domain of  $\mathbb{R}^N$ , ( $N \geq 1$ ), let  $B(x)$  be a vectorfield on  $\Omega$  whose regularity will be discussed below. We want to study the flow associated to the following ordinary differential equation

$$\frac{dX}{dt} = B(X) \quad \text{for } t \in \mathbb{R}, \quad X(0) = x \in \overline{\Omega}. \quad (31)$$

We consider here only a homogeneous equation to simplify the presentation. This is also why we assume that  $\overline{\Omega}$  is an invariant region, i.e.

$$B(x) \cdot v(x) = 0 \quad \text{on } \partial\Omega \quad (32)$$

where  $v$  is the unit exterior normal on  $\partial\Omega$ . Of course,  $X$  is a function of  $t$  and  $x \in \overline{\Omega}$  and we expect  $X(t, x)$  to belong to  $\overline{\Omega}$  for all  $t \in \mathbb{R}$ ,  $x \in \overline{\Omega}$ . We shall say that  $X$  is an a.e. flow associated to (31) if  $X \in C(\mathbb{R}; L^1(\Omega))$  and if  $X$  satisfies

$$X(t, x) \in \overline{\Omega} \quad \text{a.e. } x \in \Omega, \forall t \in \mathbb{R} \quad (33)$$

$$X(t + s, x) = X(t, X(s, x)) \quad \text{a.e. } x \in \Omega, \forall t, s \in \mathbb{R} \quad (34)$$

$$\lambda(\{x \in \Omega / X(t, x) \in N\}) = 0 \quad \text{if } \lambda(N) = 0, \forall t \in \mathbb{R} \quad (35)$$

$$\frac{\partial X}{\partial t} = B(X) \quad \text{in } \mathcal{D}'(\mathbb{R} \times \Omega) \quad (36)$$

where  $\lambda$  denotes the Lebesgue measure restricted to  $\Omega$ .

We also denote by  $Y \circ \lambda$  the image measure of  $\lambda$  by an application  $Y$  from  $\overline{\Omega}$  into  $\overline{\Omega}$  i.e. the measure on  $\overline{\Omega}$  defined by

$$\forall \varphi \in C(\overline{\Omega}), \quad \int \varphi d(Y \circ \lambda) = \int \varphi(Y(x)) dx. \quad (37)$$

We also introduce the following condition

$$X_t \circ \lambda \leq e^{C_0|t|} \lambda, \quad \forall t \in \mathbb{R} \quad (38)$$

where  $C_0 \geq 0$  and  $X_t(x) = X(t, x)$ . Of course, (38) implies (35). Then, we have the

**Theorem 4.** *Let  $B \in W^{1,1}(\Omega)$ . We assume  $\operatorname{div} B \in L^\infty(\Omega)$  and (32). We denote by  $C_0 = \|\operatorname{div} B\|_{L^\infty}$ .*

- 1) *There exists a unique a.e. flow associated to (31) satisfying (38).*
- 2) *In addition,  $X \in C^1(\mathbb{R}; L^{\frac{N}{N-1}}(\Omega)) \cap L^{\frac{N}{N-1}}(\Omega; W_{\text{loc}}^{1,\frac{N}{N-1}}(\mathbb{R}))$  and, for almost all  $x \in \Omega$ ,  $X \in C^1(\mathbb{R})$ ,  $B(X) \in C(\mathbb{R})$ , (31), (33) and (34) hold.*
- 3) *The a.e. flow  $X$  satisfies also*

$$\frac{\partial X}{\partial t} = \operatorname{div}_x(BX) - (\operatorname{div} B)X \quad \text{in } \mathcal{D}'(\mathbb{R} \times \Omega), \quad X|_{t=0} = x \quad \text{a.e. in } \Omega \quad (39)$$

and for all  $u_0 \in L^1(\Omega)$ ,  $u(X(t, x))$  is the unique renormalized solution in  $C(\mathbb{R}; L^1(\Omega))$  of

$$\frac{\partial u}{\partial t} = \operatorname{div}_x(Bu) - (\operatorname{div} B)u, \quad u|_{t=0} = u_0 \quad \text{a.e. in } \Omega. \quad (40)$$

If  $u^0 \in L^\infty(\Omega)$ ,  $u(X(t, x))$  is also the unique solution of (40) in  $L^\infty(\mathbb{R} \times \Omega)$  of (40) in the sense of distributions.

4) Let  $B_n \in W^{1,1}(\Omega)$  satisfy (32) and  $\operatorname{div} B_n \in L^\infty(\Omega)$  for all  $n \geq 1$ . We assume that  $B_n$  converges to  $B$  in  $L^1(\Omega)$  and that  $\operatorname{div} B_n$  converges to  $\operatorname{div} B$  in  $L^1(\Omega)$ . We denote by  $X_n$  the a.e. flow satisfying (38) corresponding to  $B_n$ . Then,  $X_n$  converges to  $X$  in  $C([-T, +T]; L^q(\Omega))$ , ( $\forall q < \infty$ );  $\frac{\partial X_n}{\partial t}$  converges to  $\frac{\partial X}{\partial t}$  in  $C([-T, +T]; L^1(\Omega))$  for all  $T < \infty$  and  $X_n(\cdot, x)$  converges to  $X(\cdot, x)$  uniformly on compact sets of  $\mathbb{R}$  for almost all  $x \in \Omega$ . Finally, if  $u_n^0$  converges in  $L^p(\Omega)$  to  $u^0$  ( $1 \leq p < \infty$ ), then  $u_n(t, x) = u_n^0(X_n(t, x))$  converges to  $u(t, x) = u^0(X(t, x))$  in  $C([-T, +T]; L^p(\Omega))$ , ( $\forall T < \infty$ ).

**Remarks.** 1) Analogous results hold for time-dependent vector fields  $B(t, x) \in L^1((-T, +T); W^{1,1}(\Omega))$  with  $\operatorname{div} B \in L^1((-T, +T); L^\infty(\Omega))$ , ( $\forall T < \infty$ ) or for flows in the whole space.

2) In [20], we show by a counterexample that the  $W^{1,1}(\Omega)$  regularity is in general optimal.

The very definition of renormalized solutions of (11)–(12) requires to be able to build an a.e. flow as above for the vector field  $B(t, x, v) = (v, -\nabla_x \Phi(x, t))$ . Of course, we have  $\operatorname{div}_{(x,v)} B = 0$ . Therefore, in order to prove Theorem 1 and in view of the sharpness of the  $W^{1,1}$  regularity, we have to show that  $B$  is in (say)  $L^1_{\text{loc}}(W_{\text{loc}}^{1,1})$  or in other words that  $\Phi \in L^1_{\text{loc}}(W_{\text{loc}}^{2,1})$ . But, this is not clear in

view of the estimates on  $f$  and  $\varrho$  which are basically  $L^1$  estimates. However, classical results from Harmonic analysis indicate that a rather sharp condition for this regularity to hold is to check that  $\varrho \in L^1_{\text{loc}}(\mathcal{H}^1_{\text{loc}})$  where  $\mathcal{H}^1$  denotes the multi-dimensional Hardy space. But, since  $f$  and  $\varrho$  are nonnegative, this is known to be equivalent to  $\varrho \in L^1_{\text{loc}}(L^1 \log L^1_{\text{loc}})$  where we denote by  $L^1 \log L^1_{\text{loc}}$  the set of functions  $\varphi$  such that  $|\varphi| \log |\varphi|$  is in  $L^1_{\text{loc}}$ . Again this integrability of  $\varrho$  is not clear and requires some careful analysis. In fact, this integrability of  $\varrho$  is a consequence of the estimates (24) and again this is a sharp statement making thus a rather remarkable chain of sharp statements in order to produce the desired a.e. flow (that is the particle paths). Indeed, one can show that if  $f|v|^2 \in L^1_{x,v}$ ,  $f \log f \in L^1_{x,v}$  then  $\varrho \log \varrho \in L^1_x$ , a striking fact in view of the following classical interpolation result: if  $f|v|^2 \in L^1_{x,v}$ ,  $f \in L^p_{x,v}$ ,  $(1 \leq p \leq \infty)$  then  $\varrho \in L^q_{x,v}$  with  $q = \frac{Np+2p}{Np+2}$ . Observe that  $q < p$  if  $p > 1$  and that, in some sense, this loss of integrability vanishes “before  $L^1$  at the  $L^1 \log L^1$  level”.

The final tool that we use in the proof of Theorem 1 is the improved regularity of macroscopic quantities that is velocity averages. The first indications of such a phenomenon were some results by V.I. Agoshkov [1] and by F. Golse, B. Perthame and R. Sentis [26]. The first general and sharp results were obtained in F. Golse, P.L. Lions, B. Perthame and R. Sentis [25] and partially extended in R.J. DiPerna and P.L. Lions [17], P. Gérard [24]. A complete theory is now available (R.J. DiPerna, P.L. Lions and Y. Meyer [22]) and we now present one example of such results.

**Theorem 5.** Let  $1 < p \leq 2$ , let  $\tau \in [0, 1)$ , let  $m \geq 0$ , let  $\psi \in C_0^\infty(\mathbb{R}^N)$  and let  $f(x, v)$  (resp.  $f(x, v, t) \in L^p(\mathbb{R}_x^N \times \mathbb{R}_v^N)$ ) (resp.  $L^p(\mathbb{R}_x^N \times \mathbb{R}_v^N \times \mathbb{R}_t)$ ) satisfy

$$(-\Delta_x + 1)^{-\tau/2} (-\Delta_v + 1)^{-m/2} \{v \cdot \nabla_x f\} \in L^p_{x,v} \quad (41)$$

(resp.

$$(-\Delta_{x,t} + 1)^{-\tau/2} (-\Delta_v + 1)^{-m/2} \left\{ \frac{\partial f}{\partial t} + v \cdot \nabla_x f \right\} \in L^p_{x,v,t} . \quad (42)$$

Then,  $\int_{\mathbb{R}^N} f(x, v) \psi(v) dv$  (resp.  $\int_{\mathbb{R}^N} f(x, v, t) \psi(v) dv \in B_2^{s,p}(\mathbb{R}_x^N)$  (resp.  $B_2^{s,p}(\mathbb{R}_x^N \times \mathbb{R}_t)$ ) where  $s = (1 - \tau) \frac{1}{(1+m)} \frac{1}{p'}$ .

This result proved by careful Fourier analysis admits various variants or extensions (see [22]).

We would like to conclude by indicating that these tools and the methods of proofs we used for the construction of global weak solutions admit various applications not only to the study of kinetic models and to the two categories of problems listed in the Introduction but also to other fields of Nonlinear Partial Differential Equations. One example of this fact is a recent work by P.L. Lions, B. Perthame and E. Tadmor [34] on conservation laws.

## References

1. V.I. Agoshkov: Spaces of functions with differential-difference characteristics and smoothness of solutions of the transport equation. Sov. Math. Dokl. **29** (1984) 662–666
2. L. Arkeryd: On the long time behaviour of the Boltzmann equation in a periodic box. Preprint
3. L. Arkeryd, C. Cercignani: On the convergence of solutions of the Enskog equation to solutions of the Boltzmann equation. Preprint
4. A.A. Arsenev: Global existence of a weak solution of Vlasov's systems of equations. USSR Comp. Math. Math. Phys. **15** (1975) 131–143
5. C. Bardos, F. Golse, D. Levermore: Work in preparation
6. L. Boltzmann: Weitere Studien über das Wärmegleichgewicht unter Gasmolekülen. Sitzungsberichte der Akademie der Wissenschaften Wien **66** (1872) 275–370 [Translation: Further studies on the thermal equilibrium of gas molecules. In: Kinetic Theory 2, ed. S.G. Brush. Pergamon Press, Oxford 1966]
7. T. Carleman: Problèmes mathématiques dans la théorie cinétique des gaz. Notes written by L. Carleson and O. Frostman. Publications Mathématiques de l'Institut Mittag-Leffler. Almqvist and Wiksell, Uppsala 1957
8. C. Cercignani: The Boltzmann equation and its applications. Springer, Berlin Heidelberg New York 1988
9. C. Cercignani: Work in preparation
10. S. Chapman, T.G. Cowling: The mathematical theory of non-uniform gases. 2nd edn. Cambridge Univ. Press, Cambridge 1952
11. L. Desvillettes: Convergence to equilibrium in large time for Boltzmann and B.G.K. equations. Arch. Rat. Mech. Anal. **110** (1990) 73–91
12. L. Desvillettes: On the convergence of splitting algorithms for some kinetic equations. Preprint
13. J.M. Dolbeault: Work in preparation
14. R.J. DiPerna, P.L. Lions: On the Fokker-Planck-Boltzmann equation. Comm. Math. Phys. **120** (1988) 1–23
15. R.J. DiPerna, P.L. Lions: On the Cauchy problem for Boltzmann equations: Global existence and weak stability. C.R. Acad. Sci. Paris **306** (1988) 343–346; Ann. Math. **130** (1989) 321–366
16. R.J. DiPerna, P.L. Lions: Global solutions of Boltzmann equations and the entropy inequality. Arch. Rat. Mech. Anal. (1990)
17. R.J. DiPerna, P.L. Lions: Global weak solutions of Vlasov-Maxwell systems. Comm. Pure Appl. Math. **XLII** (1989) 729–757
18. R.J. DiPerna, P.L. Lions: Solutions globales d'équations du type Vlasov-Poisson. C.R. Acad. Sci. Paris **307** (1988) 655–658; and work in preparation
19. R.J. DiPerna, P.L. Lions: Global weak solutions of kinetic equations. Sem. Matematico Torino (1990)
20. R.J. DiPerna, P.L. Lions: Ordinary differential equations, Sobolev spaces and transport theory. Invent. math. **98** (1989) 511–547
21. R.J. DiPerna, P.L. Lions: Equations différentielles ordinaires et équations de transport avec des coefficients irréguliers. In: Séminaire EDP 1988–1989, Ecole Polytechnique, Palaiseau 1989
22. R.J. DiPerna, P.L. Lions, Y. Meyer:  $L^p$  regularity of velocity averages. Ann. I.H.P. Anal. Non-Lin. (1991)
23. M.J. Esteban, B. Perthame: On the modified Enskog equation with elastic or inelastic collisions; Models with spin. Ann. I.H.P. Anal. Non-Lin. (1991)

24. P. Gérard: Moyennes de solutions d'équations aux dérivées partielles. In: Séminaire EDP 1986–1987, Ecole Polytechnique, Palaiseau 1987
25. F. Golse, P.L. Lions, B. Perthame, R. Sentis: Regularity of the moments of the solution of a transport equation. *J. Funct. Anal.* **76** (1988) 110–125
26. F. Golse, B. Perthame, R. Sentis: Un résultat de compacité pour les équations de transport et applications au calcul de la limite de la valeur propre principale d'un opérateur de transport. *C.R. Acad. Sci. Paris* **301** (1985) 341–344
27. H. Grad: Principles of the kinetic theory of gases. In: Flügge's Handbuch der Physik XII. Springer, Berlin Heidelberg New York 1958
28. K. Hamdache: Existence in the large and asymptotic behaviour for the Boltzmann equation. *Japan J. Appl. Math.* **2** (1985) 1–15
29. K. Hamdache: Global existence for weak solutions for the initial boundary value problems of Boltzmann equation. *Arch. Rat. Mech. Anal.* (1991)
30. R. Illner, M. Shinbrot: The Boltzmann equation. Global existence for a rare gas in an infinite vacuum. *Comm. Math. Phys.* **95** (1984) 217–226
31. J. Leray: Etude de diverses équations intégrales nonlinéaires et de quelques problèmes que pose l'hydrodynamique. *J. Math. Pures Appl.* **12** (1933) 1–82
32. P.L. Lions, F. Murat: Work in preparation
33. P.L. Lions, B. Perthame: Régularité des solutions du système de Vlasov-Poisson en dimension 3. *C.R. Acad. Sci. Paris*, **311** (1990) 205–210; and work in preparation
34. P.L. Lions, B. Perthame, E. Tadmor: Formulation cinétique des lois de conservation scalaires multidimensionnelles. *C.R. Acad. Sci. Paris* 1990; and work in preparation
35. J.C. Maxwell: On the dynamical theory of gases. *Phil. Trans. Roy. Soc. London* **157** (1966) 49–88
36. T. Nishida, K. Imai: Global solutions to the initial-value problem for the nonlinear Boltzmann equation. *Publ. R.I.M.S. Kyoto Univ.* **12** (1976) 229–239
37. B. Perthame: Global existence of solutions to the BGK model of Boltzmann equations. *J. Diff. Eq.* **82** (1989) 191–205
38. K. Pfaffelmoser: Globale klassische Lösungen des dreidimensionalen Vlasov-Poisson systems. Preprint
39. J. Polewczak: Global existence in  $L^1$  for the modified nonlinear Enskog equation in  $\mathbb{R}^3$ . Preprint
40. C. Truesdell, R.G. Muncaster: Fundamentals of Maxwell's kinetic theory of a simple monoatomic gas. Academic Press, New York 1980
41. S. Ukai: Solutions of the Boltzmann equation. In: Patterns and Waves. Qualitative analysis of differential equations. North-Holland, Amsterdam 1986, pp. 37–96



# Sheaf Theory for Partial Differential Equations

Pierre Schapira

Département de Mathématiques, Université Paris-Nord, F-93430 Villetaneuse, France

## 0. Introduction

In order to analyze the singularities of hyperfunction solutions of systems of partial differential equations, M. Sato introduced in 1969 the microlocalization functor and, more fundamentally, the microlocal point of view. Then began an intense activity, in what is now called “microlocal analysis”, and in the field of analytical partial differential equations the main tools, microdifferential operators and quantized contact transformations, were developed in Sato-Kawai-Kashiwara’s paper [S-K-K]. However, after their study of micro-hyperbolic systems [K-S1], M. Kashiwara and the author realized that for many problems these analytical tools were not really necessary on the condition to work with the complex of holomorphic solutions of the system,  $R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{O}_X)$ , and only to keep in mind the codirections of non-propagation of this complex, here the characteristic variety of  $\mathcal{M}$ . In other words one simply works with a complex of sheaves  $F$  on a real manifold  $X$ , and what we defined as its micro-support,  $SS(F)$ , a closed conic involutive subset of  $T^*X$ .

This was the starting point of the “microlocal study of sheaves”, developed in [K-S3, K-S4].

It is not our purpose to discuss this theory here, but we need to recall a few basic facts in order to introduce the new notion of an elliptic pair (obtained in collaboration with J.-P. Schneiders), a generalization of that of an elliptic system, the real manifold  $M$  on which the system is elliptic being replaced by an  $\mathbf{R}$ -constructible sheaf. We construct a characteristic class associated with an elliptic pair, and prove that when the pair has compact support, the complex of its holomorphic solutions has finite dimensional cohomology and the index is calculated as the integral of the characteristic class.

## 1. Microlocal Study of Sheaves

In this section we fix some notations and recall a few results of [K-S3, K-S4].

Let  $X$  be a real  $C^\infty$ -manifold. One denotes by  $\tau : TX \rightarrow X$  and by  $\pi : T^*X \rightarrow X$  its tangent and cotangent bundles, respectively. If  $M$  is a submanifold,

one denotes by  $T_M X$  and  $T_M^* X$  the normal and conormal bundles to  $M$  in  $X$ , respectively. In particular  $T_X^* X$  is the zero-section of  $T^* X$ .

One denotes by  $\Delta : \Delta \hookrightarrow X \times X$  the diagonal embedding and we identify  $X$  with  $\Delta$  and  $T^* X$  with  $T_{\Delta}^*(X \times X)$  by the first projection defined on  $X \times X$  and on  $T^*(X \times X) \simeq T^* X \times T^* X$ , respectively. If  $\Lambda$  is a subset of  $T^* X$ ,  $\Lambda^a$  will denote its image by the antipodal map on  $T^* X$ .

If  $X$  and  $Y$  are two manifolds, one denotes by  $q_1$  and  $q_2$  the first and second projection, defined on  $X \times Y$ .

Let  $A$  be a commutative unitary ring with finite global homological dimension (e.g.  $A = \mathbf{Z}$ ). One denotes by  $D(X)$  the derived category of the category of sheaves of  $A$ -modules on  $X$ , and by  $D^b(X)$  the full subcategory consisting of objects with bounded cohomology. If  $Z$  is a locally closed subset of  $X$ , one denotes by  $A_Z$  the sheaf on  $X$  which is constant with stalk  $A$  on  $Z$  and zero on  $X \setminus Z$ . One denotes by  $or_X$  the orientation sheaf on  $X$ , and by  $\omega_X$  the dualizing complex on  $X$ . Hence:

$$\omega_X \simeq or_X[\dim X],$$

where  $\dim X$  is the dimension of  $X$ .

The “six operations” (as says Grothendieck), that is, the operations  $\otimes^L$ ,  $R\mathcal{H}\text{om}$ ,  $Rf_*$ ,  $Rf_!$ ,  $f^{-1}$ ,  $f^!$ , are now classical tools that we shall not recall. We simply introduce some notations. For  $F \in \text{Ob}(D^b(X))$  and  $G \in \text{Ob}(D^b(Y))$ , one sets:

$$F \boxtimes^L G = q_1^{-1} F \otimes^L q_2^{-1} G,$$

$$D'F = R\mathcal{H}\text{om}(F, A_X),$$

$$DF = R\mathcal{H}\text{om}(F, \omega_X).$$

There are other operations of interest on sheaves. If  $M$  is a closed submanifold of  $X$  and  $F \in \text{Ob}(D^b(X))$ , the specialization of  $F$  along  $M$ ,  $v_M(F)$ , is an object of  $D^b(T_M X)$  and the microlocalization of  $F$  along  $M$ ,  $\mu_M(F)$ , an object of  $D^b(T_M^* X)$ . Sato’s functor  $\mu_M$  has been generalized in [K-S3] as follows. For  $F$  and  $G$  in  $D^b(X)$  on sets:

$$\mu \text{hom}(G, F) = \mu_A R\mathcal{H}\text{om}(q_2^{-1} G, q_1^! F).$$

Then:

$$R\pi_* \mu \text{hom}(G, F) \simeq R\mathcal{H}\text{om}(G, F),$$

$$\mu \text{hom}(A_M, F) \simeq \mu_M(F).$$

After the introduction of the functor  $\mu_M$  it became natural to work in  $T^* X$ , and M. Kashiwara and the author introduced in 1982 (cf. [K-S2]) the micro-support  $SS(F)$  of an object  $F$  of  $D^b(X)$ . This is a closed conic subset of  $T^* X$  which roughly speaking describes the set of codirections of non-propagation of  $F$ . More precisely:

**Definition 1.1.** We say that an open subset  $U$  of  $T^* X$  does not meet  $SS(F)$  if for any real  $C^1$ -function  $\varphi$  on  $X$  and any  $x_0 \in X$  such that  $(x_0; d\varphi(x_0)) \in U$ , one has:

$$\left( R\Gamma_{\{x; \varphi(x) \geq \varphi(x_0)\}}(F) \right)_{x_0} = 0.$$

An important property of the micro-support is given by:

**Theorem 1.2.** *Let  $F \in \text{Ob}(D^b(X))$ . Then  $\text{SS}(F)$  is an involutive subset of  $T^*X$ .*

(For the precise definition of “involutive”, cf. [K-S4, Ch. VI].)

One can evaluate the micro-support of sheaves after the main operations described above. For example one proves that for  $F$  and  $G$  in  $D^b(X)$ :

$$\text{SS}(\mu \hom(G, F)) \subset C(\text{SS}(F), \text{SS}(G)), \quad (1.1)$$

where  $C(\Lambda_1, \Lambda_2)$  is the normal cone of  $\Lambda_1$  along  $\Lambda_2$ , a closed subset of  $TT^*X$  that we identify with a subset of  $T^*T^*X$  by the Hamiltonian isomorphism. In particular:

$$\text{supp}(\mu \hom(G, F)) \subset \text{SS}(G) \cap \text{SS}(F). \quad (1.2)$$

Let  $f : Y \rightarrow X$  be a morphism of manifolds. One associates the maps:

$$T^*Y \xleftarrow{\lrcorner f'} Y \times_X T^*X \xrightarrow{f_\pi} T^*X \quad (1.3)$$

and one sets:

$$T_Y^*X = {}^t f'^{-1}(T^*_Y Y). \quad (1.4)$$

Using (1.1), one can evaluate the micro-support of  $f^{-1}F$  or  $f^!F$  (cf. [K-S3]). In particular if  $f$  is non-characteristic for  $F$ , that is

$$T_Y^*X \cap f_\pi^{-1}(\text{SS}(F)) \subset Y \times_X T_X^*X, \quad (1.5)$$

one gets:

$$\text{SS}(f^{-1}F) \subset {}^t f'^{-1}(\text{SS}(F)). \quad (1.6)$$

Similarly, if  $G \in \text{Ob}(D^b(Y))$  and  $f$  is proper on  $\text{supp}(G)$ , one proves:

$$\text{SS}(Rf_*G) \subset f_\pi {}^t f'^{-1}(\text{SS}(G)). \quad (1.7)$$

Remark that formulas (1.6) and (1.7) are similar to classical formulas obtained when calculating the wave front set of distributions or hyperfunctions or when calculating the characteristic variety of  $\mathcal{D}$ -modules, after non characteristic inverse images of proper direct images.

## 2. Constructible Sheaves (cf. [K-S3, K-S4])

In this section we assume all manifolds are real analytic and the base ring  $A$  is noetherian. An object  $F$  of  $D^b(X)$  is called weakly  $\mathbf{R}$ -constructible(w- $\mathbf{R}$ -

constructible for short) if there exists a subanalytic stratification  $X = \sqcup_\alpha X_\alpha$  such that for all  $\alpha$ , all  $j \in \mathbf{Z}$ , the sheaves  $H^j(F)|_{X_\alpha}$  are locally constant. If moreover for each  $x \in X$ , each  $j \in \mathbf{Z}$ , the stalk  $H^j(F)_x$  is finitely generated, one says  $F$  is  $\mathbf{R}$ -constructible. One denotes by  $D_{w-\mathbf{R}-c}^b(X)$  (resp.  $D_{\mathbf{R}-c}^b(X)$ ) the full subcategory of  $D^b(X)$  consisting of  $w$ - $\mathbf{R}$ -constructible (resp.  $\mathbf{R}$ -constructible) objects.

The involutivity Theorem 1.2 allows us to characterize microlocally  $w$ - $\mathbf{R}$ -constructible objects.

**Theorem 2.1.** *Let  $F \in \text{Ob}(D^b(X))$ . The following conditions are equivalent.*

- (a)  $F$  is  $w$ - $\mathbf{R}$ -constructible.
- (b)  $SS(F)$  is contained in a closed conic subanalytic isotropic subset of  $T^*X$ .
- (c)  $SS(F)$  is a closed conic subanalytic Lagrangian subset of  $T^*X$ .

By this result one proves easily that the category of  $w$ - $\mathbf{R}$ -constructible (resp.  $\mathbf{R}$ -constructible) sheaves is stable by the main operations on sheaves ( $Rf_*$  when  $f$  is proper,  $f^{-1}, f^!, \otimes^L, R\mathcal{H}\mathcal{O}\mathcal{M}, \mu_M, v_M, \mu \hom$ ).

If  $X$  is a complex manifold one defines similarly the notions of  $w$ - $\mathbf{C}$ -constructible and  $\mathbf{C}$ -constructible sheaves, by assuming that the stratas of the subanalytic stratification  $X = \sqcup_\alpha X_\alpha$  are complex analytic submanifolds. Then the link between  $\mathbf{R}$ - and  $\mathbf{C}$ -constructibility is given by:

**Theorem 2.2.** *Let  $F \in \text{Ob}(D_{w-\mathbf{R}-c}^b(X))$ . Then  $F$  is  $w$ - $\mathbf{C}$ -constructible if and only if  $SS(F)$  is conic for the action of  $\mathbf{C}^\times$  on  $T^*X$ .*

**Remark 2.3.** Note that the microlocal study of constructible sheaves was initiated by Kashiwara [K1].

### 3. $\mathcal{D}$ -Modules

We shall not review this theory here and refer to [S-K-K, K1, S1] for detailed expositions. We shall only fix a few notations and make the link with the micro-support.

From now on the base ring  $A$  is the field  $\mathbf{C}$  of complex numbers.

Let  $(X, \mathcal{O}_X)$  be a complex manifold of complex dimension  $n$ . One denotes by  $\mathcal{D}_X$  (resp.  $\mathcal{D}_X^{op}$ ) the sheaf on  $X$  of finite order (resp. infinite order) holomorphic differential operators and one sets  $\Omega_X = \mathcal{O}_X^{(n)} \otimes \mathcal{O}_X$ , where  $\mathcal{O}_X^{(n)}$  is the sheaf of holomorphic  $n$ -forms. One denotes by  $D(\mathcal{D}_X)$  (resp.  $D(\mathcal{D}_X^{op})$ ) the derived category of the abelian category of left (resp. right)  $\mathcal{D}_X$ -modules, and by  $D^b(\mathcal{D}_X)$  (resp.  $D_{coh}^b(\mathcal{D}_X)$ ) the full triangulated subcategory of  $D(\mathcal{D}_X)$  consisting of objects with bounded (resp. bounded and coherent) cohomology. One defines similarly  $D^b(\mathcal{D}_X^{op})$  and  $D_{coh}^b(\mathcal{D}_X^{op})$ .

If  $\mathcal{M}$  is an object of  $D_{coh}^b(\mathcal{D}_X)$ , its characteristic variety denoted  $\text{char}(\mathcal{M})$ , is a closed conic analytic subset of  $T^*X$ , which is involutive ([S-K-K]). In fact one has:

**Theorem 3.1** ([K-S2]). Let  $\mathcal{M} \in \text{Ob}(D_{\text{coh}}^b(\mathcal{D}_X))$ . Then:

$$\text{SS}(R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{O}_X)) = \text{char}(\mathcal{M}).$$

Note that the inclusion  $\subset$  in Theorem 3.1, which is the most useful for applications, is easily deduced from the Cauchy-Kowalevski theorem, in its precised form due to Leray [L]. The converse inclusion makes use of the sheaf of rings  $\mathcal{E}_X^R$  of [S-K-K]. Also note that this result, combined with Theorem 1.2, gives a new proof of the involutivity of the characteristic variety of  $\mathcal{D}$ -modules.

Let  $f : Y \rightarrow X$  be a morphism of complex manifolds. One denotes by  $\mathcal{D}_{Y \rightarrow X}$  the sheaf  $\mathcal{O}_Y \otimes_{f^{-1}\mathcal{O}_X} f^{-1}\mathcal{D}_X$  endowed with its natural structure of a  $(\mathcal{D}_Y, f^{-1}\mathcal{D}_X)$ -bimodule.

Let  $\mathcal{M} \in \text{Ob}(D^b(\mathcal{D}_X))$ . One sets:

$$\underline{f^{-1}\mathcal{M}} = \mathcal{D}_{Y \rightarrow X} \otimes_{f^{-1}\mathcal{D}_X}^L f^{-1}\mathcal{M}.$$

Let  $\mathcal{N} \in \text{Ob}(D^b(\mathcal{D}_Y^{op}))$ . One sets

$$\underline{f_*\mathcal{N}} = Rf_*(\mathcal{N} \otimes_{\mathcal{D}_Y}^L \mathcal{D}_{Y \rightarrow X}).$$

If  $\mathcal{M} \in \text{Ob}(D^b(\mathcal{D}_X))$  and  $\mathcal{N} \in \text{Ob}(D^b(\mathcal{D}_Y))$  one sets:

$$\underline{\mathcal{M} \boxtimes \mathcal{N}} = \mathcal{D}_{X \times Y} \otimes_{\mathcal{D}_X \boxtimes \mathcal{D}_Y} (\mathcal{M} \boxtimes \mathcal{N}),$$

and there is a similar formula for right modules.

Finally one sets:

$$\begin{aligned} \underline{D'\mathcal{M}} &= R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{D}_X), \\ \underline{D\mathcal{M}} &= R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \Omega_X \otimes_{\mathcal{O}_X} \mathcal{D}_X[n]). \end{aligned}$$

(In this last formula,  $\mathcal{M}$  is a right module.)

#### 4. Microfunction Solutions of $\mathcal{D}$ -Modules

Let  $M$  be a real analytic manifold,  $X$  a complexification of  $M$ ,  $\mathcal{M}$  a left coherent  $\mathcal{D}_X$ -module. By considering the complex  $R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{O}_X)$  and using Theorem 3.1, one may recover many classical results. For example, applying (1.2) with  $G = D'(C_M)$  one gets:

$$\text{supp}(R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{C}_M)) \subset T_M^*X \cap \text{char}(\mathcal{M}),$$

where  $\mathcal{C}_M$  is the sheaf of Sato microfunctions. In particular this shows that the analytic wave front set of a hyperfunction solution of a system of linear differential equations is contained in the characteristic variety of the system (“Sato’s principle”). More generally, the inclusion (1.1) immediately implies that the microfunction solutions of the system  $\mathcal{M}$  extend in the micro-hyperbolic directions, and one recovers the results of [K-S1] in the differential case. Microdifferential

systems can be treated similarly, once the microlocal action of  $\mathcal{E}_X^R$  on  $\mathcal{O}_X$  is defined as in [K-S3].

These techniques can also be applied to the study of boundary value problems and diffraction, including the case of non smooth obstacles. It is then useful to introduce new sheaves of microfunctions (using the functor  $\mu \text{hom}$ ) and new wave front sets. We refer to [S2] for details.

## 5. Elliptic Pairs

In this section we expose new results obtained in collaboration with J.-P. Schneiders. Let  $X$  be a complex manifold of complex dimension  $n$ . If there is no risk of confusion, we identify  $X$  with the real underlying manifold.

**Definition 5.1.** An elliptic pair  $(\mathcal{M}, F)$  on  $X$  is the data of  $\mathcal{M} \in \text{Ob}(D_{\text{coh}}^b(\mathcal{D}_X))$  and  $F \in \text{Ob}(D_{R-c}^b(X))$  satisfying:  $\text{char}(\mathcal{M}) \cap SS(F) \subset T_X^*X$ .

We use the same terminology for objects of  $D_{\text{coh}}^b(\mathcal{D}_X^{op})$ .

**Theorem 5.2** (cf. [S-Sc]). *Let  $(\mathcal{M}, F)$  be an elliptic pair.*

(i) *Regularity. The natural morphism:*

$$R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M}, D'F \otimes \mathcal{O}_X) \rightarrow R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{O}_X)$$

*is an isomorphism.*

(ii) *Assume  $\text{supp}(\mathcal{M}) \cap \text{supp}(F)$  is compact.*

(a) *Finiteness. For all  $j \in \mathbb{Z}$ , the  $\mathbf{C}$ -vector spaces*

$$H^j R\Gamma(X; R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{O}_X))$$

*are finite dimensional.*

(b) *Duality. The pairing*

$$R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{O}_X) \otimes (\Omega_X \otimes_{\mathcal{D}_X}^L \mathcal{M} \otimes F) \longrightarrow \Omega_X \otimes_{\mathcal{D}_X}^L \mathcal{O}_X$$

*and the integration morphism  $H_c^n(X; \Omega_X \otimes_{\mathcal{D}_X}^L \mathcal{O}_X) \simeq H_c^{2n}(X; \text{or}_X) \rightarrow \mathbf{C}$  induce a perfect duality on the spaces  $H^j R\Gamma(X; R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{O}_X))$  and  $H^{n-j} R\Gamma(X; \Omega_X \otimes_{\mathcal{D}_X}^L \mathcal{M} \otimes F)$ .*

(c) *Parameters. Let  $Y$  be another complex manifold. Then the natural morphism*

$$(Rq_{2*}q_1^{-1} R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{O}_X)) \otimes \mathcal{O}_Y \longrightarrow Rq_{2*} R\mathcal{H}\text{om}_{q_1^{-1}\mathcal{D}_X}(q_1^{-1}(\mathcal{M} \otimes F), \mathcal{O}_{X \times Y})$$

*is an isomorphism.*

*Sketch of Proof.* (i) follows from a general result of [K-S4] which asserts that the natural morphism

$$R\mathcal{H}\text{om}(F, A_X) \otimes^L G \rightarrow R\mathcal{H}\text{om}(F, G)$$

is an isomorphism as soon as  $SS(F) \cap SS(G)$  is contained in  $T_X^*X$ .

(ii) Using techniques of [Sc] one can reduce the problem to the case where  $\mathcal{M}$  admits a free presentation. Next by adding the Cauchy-Riemann system to  $\mathcal{M}$ , one can assume  $F$  is supported by a real analytic manifold  $M$  whose complexification is  $X$ . Then one represents  $F$  by a bounded complex whose components are direct sums of sheaves  $\mathbf{C}_U$ ,  $U$  open, relatively compact, subanalytic in  $M$  and such that  $D'_M \mathbf{C}_U = \mathbf{C}_{\bar{U}}$ , (here  $D'_M$  is the duality functor on  $M$ ). Then (ii) is proved by similar arguments as those in [B-S] or [R-R].  $\square$

By adapting Kashiwara's construction of the characteristic cycle of  $\mathbf{R}$ -constructible sheaves (cf. [K-S4, Ch. IX]) we shall now construct a characteristic class associated with an elliptic pair.

Let  $(\mathcal{M}, F)$  be an elliptic pair and assume  $\mathcal{M}$  is a right module. Sato's isomorphism  $\mathcal{D}_X^\infty \simeq \delta^! \mathcal{O}_{X \times X}^{(0,n)}[n]$  induces a morphism:

$$R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{M} \otimes F) \rightarrow \delta^!(\underline{D}(\mathcal{M} \otimes F) \boxtimes (\mathcal{M} \otimes F) \otimes_{\mathcal{D}_{X \times X}}^L \mathcal{O}_{X \times X}). \quad (5.1)$$

Moreover the natural morphism:

$$\Omega_X \otimes_{\mathcal{O}_X} \mathcal{D}_X[n] \otimes_{\mathcal{D}_X \otimes \mathcal{D}_X}^L \delta^{-1} \mathcal{O}_{X \times X} \rightarrow \Omega_X \otimes_{\mathcal{D}_X}^L \mathcal{O}_X[n]$$

induces a morphism:

$$\delta^{-1}(\underline{D}(\mathcal{M} \otimes F) \boxtimes (\mathcal{M} \otimes F) \otimes_{\mathcal{D}_{X \times X}}^L \mathcal{O}_{X \times X}) \rightarrow \omega_X. \quad (5.2)$$

Set for short:

$$H = \underline{D}(\mathcal{M} \otimes F) \boxtimes (\mathcal{M} \otimes F) \otimes_{\mathcal{D}_{X \times X}}^L \mathcal{O}_{X \times X}. \quad (5.3)$$

We get the chain of morphisms:

$$\begin{aligned} R\mathcal{H}\text{om}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{M} \otimes F) &\rightarrow \delta^! H \\ &\rightarrow \delta^! \delta_* \delta^{-1} H \simeq \delta^{-1} H \\ &\rightarrow \omega_X. \end{aligned}$$

Hence setting:

$$S = \text{char}(\mathcal{M}) \cap \text{supp}(F), \quad (5.4)$$

we get a morphism:

$$\text{Hom}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{M} \otimes F) \rightarrow H_S^0(X; \omega_X). \quad (5.5)$$

In fact this construction can be made "microlocal". Set:

$$A = \text{char}(\mathcal{M}) + SS(F)^a. \quad (5.6)$$

Since the micro-support of  $H$  is contained in  $A \times A^a$ , we have the isomorphisms:

$$\begin{aligned} \delta^! H &\simeq R\pi_* \mu_A H \\ &\xleftarrow{\sim} R\pi_* R\Gamma_A \mu_A H, \end{aligned}$$

and we get the morphism:

$$\mathrm{Hom}_{\mathcal{D}_X}(\mathcal{M} \otimes F, \mathcal{M} \otimes F) \rightarrow H_A^0(T^*X; \pi^{-1}\omega_X). \quad (5.7)$$

**Definition 5.3.** We call the image of 1 in  $H_S^0(X; \omega_X)$  (resp. in  $H_A^0(T^*X; \pi^{-1}\omega_X)$ ) by the morphism (5.5) (resp. (5.7)) the Euler class (resp. the microlocal Euler class) of the elliptic pair  $(\mathcal{M}, F)$  and we denote it by  $\mathrm{eu}(\mathcal{M}, F)$  (resp.  $\mu\mathrm{eu}(\mathcal{M}, F)$ ).

Of course  $\mathrm{eu}(\mathcal{M}, F)$  is the restriction of  $\mu\mathrm{eu}(\mathcal{M}, F)$  to the zero-section of  $T^*X$ .

**Definition 5.4.** Let  $(\mathcal{M}, F)$  be an elliptic pair such that  $\mathrm{supp}(\mathcal{M}) \cap \mathrm{supp}(F)$  is compact and assume  $\mathcal{M}$  is a right module. One sets:

$$\chi(X; \mathcal{M}, F) = \sum_j (-1)^j \dim(H^j(R\Gamma(X; F \otimes \mathcal{M} \otimes_{\mathcal{D}_X}^L \mathcal{O}_X))).$$

By adapting to the case of  $\mathcal{D}$ -modules Kashiwara's proof of the index theorem for constructible sheaves (cf. [K-S4, Ch. IX]), we can prove:

**Theorem 5.5.** *In the situation of Definition 5.4, one has:*

$$\chi(X; \mathcal{M}, F) = \int_X \mathrm{eu}(\mathcal{M}, F). \quad (5.8)$$

**Examples and Comments 5.6.** (a) Assume  $X$  is the complexification of a real analytic manifold  $M$ , and let  $\mathcal{M}$  be a coherent  $\mathcal{D}_X$ -module. Then  $\mathcal{M}$  is elliptic on  $M$  in the classical sense if and only if  $(\mathcal{M}, \mathbf{C}_M)$  is an elliptic pair, which simply means that:

$$\mathrm{char}(\mathcal{M}) \cap T_M^*X \subset T_X^*X.$$

In this case the isomorphism of Theorem 5.2 applied to the elliptic pair  $(\mathcal{M}, D'\mathbf{C}_M)$  gives the isomorphism:

$$R\mathrm{Hom}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{A}_M) \simeq R\mathrm{Hom}_{\mathcal{D}_X}(\mathcal{M}, \mathcal{B}_M)$$

where  $\mathcal{A}_M$  (resp.  $\mathcal{B}_M$ ) denotes the sheaf of real analytic functions (resp. hyperfunctions) on  $M$ . If  $M$  is compact, one recovers the classical finiteness theorem for elliptic systems, and the index is calculated by the Atiyah-Singer theorem [A-S].

(b) Let  $\Omega$  be an open subset of  $X$  with real analytic boundary. Then  $(\mathcal{M}, \mathbf{C}_\Omega)$  is an elliptic pair if and only if  $\partial\Omega$  is non-characteristic for  $\mathcal{M}$ , that is:

$$\mathrm{char}(\mathcal{M}) \cap T_{\partial\Omega}^*X \subset T_X^*X.$$

If  $\Omega$  is relatively compact, one gets the finiteness of the spaces  $\mathrm{Ext}_{\mathcal{D}_X}^j(\Omega, \mathcal{M}, \mathcal{O}_X)$ , a result of Bony and Schapira [B-S] (in case  $\mathcal{M} = \mathcal{D}_X/\mathcal{D}_X.P$ , cf. Kashiwara [K1] and Kawai [Ka] for various generalizations), extended to the relative case by

Houzel and Schapira [H-S], the index being calculated by Boutet de Monvel and Malgrange [B-M].

(c) For any  $\mathcal{M} \in \text{Ob}(D_{\text{coh}}^b(\mathcal{D}_X))$ ,  $(\mathcal{M}, C_X)$  is an elliptic pair. In this case the duality theorem is due to Mebkhout [M]. If  $\mathcal{M}$  has compact support one recovers many classical results. In particular if  $\mathcal{G}$  is a coherent  $\mathcal{O}_X$ -module with compact support, one can apply the theorem with  $\mathcal{M} = \mathcal{D}_X \otimes_{\mathcal{O}_X} \mathcal{G}$  and recover theorems of Cartan and Serre (cf. [C-S, Se]). Concerning the index, let us recall that O'Brian, Toledo and Tong [O-T-T], generalizing the Hirzebruch-Riemann-Roch formula [H], constructed the Chern class of coherent  $\mathcal{O}_X$ -modules with compact support, and proved an index theorem in this case. For the case of  $\mathcal{D}_X$ -modules with compact support, cf. Angénol-Lejeune [A-L].

(d) For any  $F \in \text{Ob}(D_{R-c}^b(X))$ ,  $(\mathcal{O}_X, F)$  is an elliptic pair and its microlocal Euler class coincides with the Lagrangian cycle of  $F$  defined by Kashiwara in [K2]. Remark that if  $G$  is an  $R$ -constructible object on a real manifold  $M$ , one can associate to it an elliptic pair, namely  $(\mathcal{O}_X, i_* G)$  where  $i : M \hookrightarrow X$  is a complexification of  $M$ . In this case the index is calculated by Kashiwara (loc. cit.) (cf. also Dubson [D] and Ginsburg [G] in the complex case).

(e) Let  $B(x_0; \varepsilon)$  denote the open ball (in a local chart at  $x_0$ ) with center  $x_0$  and radius  $\varepsilon$  and let  $\mathcal{M}$  be an holonomic  $\mathcal{D}_X$ -module. Then for  $\varepsilon$  small enough,  $(\mathcal{M}, C_{B(x_0; \varepsilon)})$  is an elliptic pair (cf. [K1]).

## References

- [A-L] B. Angénol, M. Lejeune-Jalabert: Le théorème de Riemann-Roch singulier pour les  $D$ -modules. In: Systèmes différentiels et singularités. Astérisque **130** (1985) 130–160
- [A-S] M. Atiyah, I.M. Singer: The index of elliptic operators on compact manifolds. Bull. Am. Math. Soc. **69** (1963) 422–433
- [B-M] M. Boutet de Monvel, B. Malgrange: Le théorème de l'indice relatif. Ann. Sc. Ec. Norm. Sup. **23** (1990) 151–192
- [B-S] J.-M. Bony, P. Schapira: Existence et prolongement des solutions holomorphes des équations aux dérivées partielles. Invent. math. **17** (1972) 95–105
- [C-S] H. Cartan, J.-P. Serre: Un théorème de finitude concernant les variétés analytiques compactes. C. R. Acad. Sci. **237** (1953) 128–130
- [D] A. Dubson: Formules pour l'indice des faisceaux constructibles et  $\mathcal{D}$ -modules holonomes. C. R. Acad. Sci. **298** (1984) 113–116
- [G] V. Ginsburg: A theorem on the index of differential systems and the geometry of manifolds with singularities. Sov. Math. Dokl. **31** (1985) 309–313
- [H] F. Hirzebruch: Topological methods in algebraic geometry. (Grundlehren der mathematischen Wissenschaften, vol. 131.) Springer, Berlin Heidelberg New York 1966
- [K1] M. Kashiwara: Systems of microdifferential equations. Progress in Math., Birkhäuser, 1983
- [K2] M. Kashiwara: Index theorem for constructible sheaves. In: Systèmes différentiels et singularités. Astérisque **130** (1985) 193–209
- [K-S1] M. Kashiwara, P. Schapira: Micro-hyperbolic systems. Acta Math. **142** (1979) 1–55
- [K-S2] M. Kashiwara, P. Schapira: Micro-support des faisceaux. C. R. Acad. Sci. **295** (1982) 487–490

- [K-S3] M. Kashiwara, P. Schapira: Microlocal study of sheaves. Astérisque **128** (1985)
- [K-S4] M. Kashiwara, P. Schapira: Sheaves on manifolds. (Grundlehren der mathematischen Wissenschaften, vol. 292.) Springer, Berlin Heidelberg New York 1990
- [Ka] T. Kawai: Finite dimensionality of cohomology groups attached to systems of linear differential equations. J. Math. Kyoto Univ. **13** (1973) 73–95
- [L] J. Leray: Problème de Cauchy I. Bull. Soc. Math. France **85** (1957) 389–430
- [M] Z. Mebkhout: Théorèmes de dualité globale pour les  $D$ -modules cohérents. Math. Scand. **50** (1982) 25–43
- [O-T-T] N. O'Brian, D. Toledo, Y.L. Tong: Hirzebruch-Riemann-Roch for coherent sheaves. Amer. J. Math. **103** (1981) 253–271
- [R-R] J.-P. Ramis, G. Ruget: Complexe dualisant et théorèmes de dualité en géométrie analytique complexe. Publ. Math. I.H.E.S. **38** (1971) 77–91
- [S-K-K] M. Sato, T. Kawai, M. Kashiwara: Hyperfunctions and pseudo-differential equations. (Lecture Notes in Mathematics, vol. 287.) Springer, Berlin Heidelberg New York 1973
- [S1] P. Schapira: Microdifferential systems in the complex domain. (Grundlehren der mathematischen Wissenschaften, vol. 269.) Springer, Berlin Heidelberg New York 1985
- [S2] P. Schapira: Microfunctions for boundary value problems. In: Algebraic Analysis, papers dedicated to Prof. Sato (M. Kashiwara and T. Kawai, eds.). Academic Press, New York 1988, pp. 809–819
- [S-Sc] P. Schapira, J.-P. Schneiders: Paires elliptiques I – Finitude et dualité. C. R. Acad. Sci. **311** (1990) 83–86; II – Classe d'Euler et indice. C. R. Acad. Sci. **312** (1991) 81–84
- [Sc] J.-P. Schneiders: Un théorème de dualité relative pour les modules différentiels. C. R. Acad. Sci. **303** (1986) 235–238. Thèse, Univ. Liège 1986 et article à paraître.
- [Se] J.-P. Serre: Un théorème de dualité. Comm. Math. Helv. **29** (1954) 9–26

# The Evolution of Harmonic Maps

Michael Struwe

Mathematik, ETH-Zentrum, CH-8092 Zürich, Switzerland

1. Let  $M$  and  $N$  be compact Riemannian manifolds of dimensions  $m$  and  $\ell$  and with metrics  $\gamma$  and  $g$ , respectively. We may assume that  $N$  is isometrically embedded in some Euclidean  $\mathbb{R}^n$ .

For  $C^1$ -maps  $u : M \rightarrow N \subset \mathbb{R}^n$  let

$$e(u) = \frac{1}{2} |\nabla u|^2 = \frac{1}{2} \sum_{\alpha, \beta=1}^m \sum_{i=1}^n \gamma^{\alpha\beta} \frac{\partial}{\partial x_\alpha} u^i \frac{\partial}{\partial x_\beta} u^i \quad (1.1)$$

be the energy density and – with  $dM = \sqrt{\det(\gamma_{\alpha\beta})} dx$  – let

$$E(u) = \int_M e(u) dM$$

be the energy of  $u$ . If  $M \subset \mathbb{R}^m$  carries the Euclidean metric,  $E$  is nothing but the standard Dirichlet integral

$$E(u) = \frac{1}{2} \int_M |\nabla u|^2 dx.$$

A map  $u$  is harmonic if  $E$  is stationary at  $u$ , which for  $u \in C^2$  is equivalent to the condition that the vector  $(\Delta_M u)(x)$  at all points  $x \in M$  is orthogonal to the tangent space  $T_{u(x)}N$  to  $N$  at  $u(x) \in \mathbb{R}^n$ ; that is,

$$\Delta_M u \perp T_u N. \quad (1.2)$$

( $\Delta_M$  is the Laplace-Beltrami operator on  $M$ .) In local coordinates on  $N$ , (1.2) takes the form

$$\Delta_M u = \Gamma(u)(\nabla u, \nabla u) \quad (1.3)$$

with a bilinear form  $\Gamma$  involving the Christoffel symbols of the metric  $g$  on  $N$ . Again, if  $M \subset \mathbb{R}^m$  and if, in particular,  $N = S^\ell \subset \mathbb{R}^n$  ( $n = \ell + 1$ ), (1.3) takes the following simple form

$$\Delta u = u |\nabla u|^2.$$

The notion of harmonic map generalizes the concepts of closed geodesic (for  $M = S^1$ ) and harmonic function (for  $N = \mathbb{R}$ ).

The study of harmonic maps was initiated by Fuller, Nash, and Sampson; see Eells-Lemaire [12] for further background material and references. The first general existence result is due to Eells and Sampson. Since standard variational techniques fail, in their pioneering work Eells and Sampson [13] introduce the evolution problem

$$\partial_t u - \Delta_M u \perp T_u N \quad \text{on } M \times ]0, \infty[ \quad (1.4)$$

with initial condition

$$u = u_0 \quad \text{at } t = 0. \quad (1.5)$$

Upon multiplying (1.4) by  $\partial_t u \in T_u N$  and integrating by parts, we at once obtain the energy inequality

$$E(u(T)) + \int_0^T \int_M |\partial_t u|^2 dM dt \leq E(u_0), \quad \forall T > 0 \quad (1.6)$$

for any solution  $u$  of (1.4), (1.5). That is, (1.4) is the  $L^2$ -gradient flow for  $E$ . In local coordinates again, (1.4) takes the form

$$\partial_t u - \Delta_M u = \Gamma(u)(\nabla u, \nabla u). \quad (1.7)$$

Upon differentiation the latter expression and multiplying by  $\nabla u$ , we moreover obtain the following Bochner-type differential inequality for the energy density

$$(\partial_t - \Delta_M)e(u) + c|\nabla^2 u|^2 \leq K_N e(u)^2 + C_M e(u), \quad (1.8)$$

where  $K_N$  denotes an upper bound for the sectional curvature of  $N$ ,  $c > 0$ , and  $C_M$  depends on (the Ricci curvature of) the metric  $\gamma$ . (1.6), (1.8) and Moser's weak Harnack inequality for parabolic equations now lead to the following result.

**Theorem 1** (Eells-Sampson [13]). *Suppose  $K_N \leq 0$ . Then for any smooth map  $u_0 : M \rightarrow N$  problem (1.4), (1.5) admits a unique, smooth solution  $u(t)$  which, as  $t \rightarrow \infty$  suitably, converges to a smooth harmonic map  $u_\infty$  homotopic to  $u_0$ .*

Theorem 1 was extended to manifolds with boundary by Hamilton [17]. Moreover, in this case the curvature restriction  $K_N \leq 0$  can be weakened if the initial and boundary data have small range (Jost [21]). However, for general data, the curvature restriction  $K_N \leq 0$  is in some sense optimal as Eells and Wood [14] show that Theorem 1 ceases to be true for  $M = T^2$ ,  $N = S^2$  and an initial map  $u_0$  of topological degree 1.

**2.** Nevertheless, in two dimensions ( $m = 2$ ) Lemaire [23] and Sacks-Uhlenbeck [29] independently showed that also the topological condition  $\pi_2(N) = 0$  suffices to find harmonic representatives for all homotopy classes of maps  $u_0 : M \rightarrow N$ . In fact, a new proof of this result, using the evolution problem (1.4), can be given [35]. (Since by the result of Eells and Wood we must expect singularities, we consider also maps in the Sobolev space

$$H^{1,2}(M; N) = \left\{ u \in H^{1,2}(M; \mathbb{R}^n); \quad u(M) \subset N \right\}$$

of measurable, finite energy maps  $u : M \rightarrow N$ ; that is, with distributional derivative in  $L^2$ .) Then we have the following generalization of Theorem 1.

**Theorem 2** ([33; Theorem 4.2]). Suppose  $m = 2$ . For any  $u_0 \in H^{1,2}(M; N)$  there exists a global weak solution  $u : M \times [0, \infty) \rightarrow N$  of (1.4), (1.5) which satisfies (1.6) and is  $(C^\infty)$ -regular away from finitely many points  $(\bar{x}_k, \bar{t}_k)$ ,  $1 \leq k \leq K$ . The solution  $u$  is unique in this class.

At a singularity  $(\bar{x}, \bar{t})$  a “harmonic sphere”  $\bar{u} : S^2 \cong \overline{\mathbb{R}^2} \rightarrow N$  separates in the sense that for suitable  $x_m \rightarrow \bar{x}$ ,  $R_m \searrow 0$ ,  $t_m \nearrow \bar{t}$  we have

$$u_m(x) := u(\exp_{x_m}(R_m x), t_m) \longrightarrow \bar{u} \quad \text{in } H_{\text{loc}}^{2,2}(\mathbb{R}^2; N), \quad (2.1)$$

where  $\bar{u} \not\equiv \text{const.}$  is harmonic, has finite energy and extends to a smooth harmonic map  $\bar{u} : S^2 \cong \overline{\mathbb{R}^2} \rightarrow N$ .

Finally, as  $t \rightarrow \infty$  suitably,  $u(t) \rightarrow u_\infty$  weakly in  $H^{1,2}(M; N)$ , where  $u_\infty : M \rightarrow N$  is smooth and harmonic. Convergence is strong away from finitely many points  $(\bar{x}_\ell, \bar{t}_\ell = \infty)$ ,  $1 \leq \ell \leq L$ , where harmonic spheres separate in the sense (2.1).

**Remark.** Let  $\varepsilon_0 = \inf\{E(u); u : S^2 \rightarrow N \text{ is non-constant and harmonic}\} > 0$ . Then by (2.1) we can bound  $K + L \leq \varepsilon_0^{-1} E(u_0)$ . In particular, for initial data such that  $E(u_0) < \varepsilon_0$  the solution  $u$  constructed in Theorem 2 is smooth and converges uniformly to a harmonic limit.

Theorem 2 was extended to 2-manifolds  $M$  with boundary  $\partial M \neq \emptyset$  by Chang [3] with applications to results by Brezis-Coron and Jost on the Dirichlet problem for harmonic maps into the sphere, and by Chen-Musina [7] to the case of target manifolds with boundary.

Moreover, variants of (1.4) arise if one attempts to solve free boundary problems for minimal surfaces by a deformation method and results completely analogous to Theorem 2 hold; see [34]. M. Li [24] has recently generalized these results to free boundary problems for harmonic maps.

To this day it is not known whether in general (in two dimensions) the flow (1.4) will encounter singularities in finite time. (Of course, the result of Eells and Wood provides us with an example where either such singularities exist or the flow fails to converge asymptotically.) However, some recent results of Chang and Ding [4], respectively work by Grayson and Hamilton [16] lend support to the conjecture that for  $m = 2$  the flow (1.4) does not develop singularities in finite time.

3. The situation is quite different in higher dimensions, as there are energy minimizing harmonic maps with singularities. An example is given by the well-known map  $u(x) = \frac{x}{|x|} : B_1(0) \subset \mathbb{R}^m \rightarrow S^{m-1} \subset \mathbb{R}^m$  (Brezis-Coron-Lieb [2], Lin [26]); if  $m \geq 7$ ,  $u$  is minimizing even if we regard  $u$  as a map  $u : B_1(0) \rightarrow S^m \subset \mathbb{R}^{m+1}$  (Jäger-Kaul [20]). For energy-minimizing maps therefore only partial regularity results can be expected. Such a regularity theory was developed by Schoen and Uhlenbeck [30] – and independently by Giaquinta-Giusti [15] in the case that the image is covered by a single chart. A crucial role is played by a subtle monotonicity estimate.

Similarly, progress on the evolution problem (1.4) for  $m \geq 3$  and general targets came as a consequence of a peculiar monotonicity formula for (1.4), discovered in [36; Lemma 3.2] for maps  $u : \mathbb{R}^m \times [0, T] \rightarrow N$  and extended to curved domains by Chen-Struwe [8]. Let

$$G_{z_0}(x, t) = \frac{1}{(4\pi(t_0 - t))^{m/2}} \exp\left(-\frac{|x - x_0|^2}{4(t_0 - t)}\right), \quad \text{if } t < t_0 \quad (3.1)$$

be the fundamental solution to the heat equation on  $\mathbb{R}^m \times \mathbb{R}$  with singularity at  $z_0 = (x_0, t_0)$ , and let  $\varphi \in C_0^\infty(\mathbb{R}^m)$  be a smooth cut-off function such that  $\varphi \equiv 1$  in a neighbourhood of 0 and such that the support of  $\varphi$  is contained in a ball of radius  $\varrho$  less than the injectivity radius  $\varrho_M$  of  $M$ . For a solution  $u : M \times [0, T] \rightarrow N$  and  $R^2 < t_0 < T$  let

$$\Phi(R; z_0) = \frac{1}{2} R^2 \int_{B_\varrho(x_0)} |\nabla u|^2 \varphi^2 G_{z_0} dx|_{t=t_0-R^2},$$

in local normal coordinates around  $x_0$  on  $M$ .

**Theorem 3** (Chen-Struwe [8; Lemma 4.2]). *There exists a constant  $C$  depending only on  $M$  and  $N$  such that for any  $T > 0$ , any  $0 < R^2 < R_0^2 \leq t_0 < T$  and any regular solution  $u : M \times [0, T] \rightarrow N$  of (1.4) there holds*

$$\Phi(R; z_0) \leq \exp(C(R_0 - R)) \Phi(R_0; z_0) + CE(u_0)(R_0 - R). \quad (3.2)$$

A particular consequence of the monotonicity formula (3.2) is the following.

**Theorem 4** ([36; Theorem 5.1]; Chen-Struwe [8; Lemma 4.4]). *There exists a constant  $\varepsilon_0 > 0$  depending only on  $M$  and  $N$  such that for any solution  $u : M \times [0, T] \rightarrow N$  of (1.4) the following is true:*

*If  $\Phi(R; z_0) < \varepsilon_0$  for some  $z_0 = (x_0, t_0)$ ,  $0 < R^2 \leq t_0 < T$ ,  $R < \varrho_M$ , then  $|\nabla u(z_0)| \leq C$ , with a constant  $C = C(M, N, R)$ .*

Theorem 4 at once leads to a new proof of Mitteau's [27] global existence result for initial data with small energy density. More important, Theorem 4 together with a penalization device to obtain approximate solutions for (1.4), due to Chen [5], Keller-Rubinstein-Sternberg [22], and Shatah [31] led to a global existence and partial regularity result for (1.4) in higher dimensions ( $m \geq 3$ ).

**Theorem 5** (Chen-Struwe [8; Theorem 1.5]). *For any (smooth) map  $u_0 : M \rightarrow N$  there exists a global weak solution  $u : M \times [0, \infty] \rightarrow N$  of (1.4), (1.5) satisfying (1.6) and regular off a set of co-dimension  $\geq 2$ . (In fact, this dimension estimate holds for all  $t > 0$ ; see Cheng [9].) As  $t \rightarrow \infty$  suitably,  $u(t) \rightarrow u_\infty$  weakly in  $H^{1,2}(M, N)$  where  $u_\infty : M \rightarrow N$  is harmonic and regular off a set of co-dimension  $\geq 2$ .*

Moreover,  $u$  satisfies a variant of the monotonicity formula (3.2). This fact was used by Coron [10] to prove that the solution obtained in Theorem 5 is in general not unique – even among partially regular solutions satisfying (1.6).

The estimate on the co-dimension of the singular set very likely can be improved to 3, as for energy-minimizing (stationary) harmonic maps; see Giacinti-Giusti [15] and Schoen-Uhlenbeck [30]. However, as was first observed by Coron-Ghidaglia [11], in higher dimensions singularities may appear in finite time. Subsequently, Chen and Ding [6] gave an argument relating singularities to the fact that in higher dimensions the infimum of the energy in certain non-trivial homotopy classes of maps may be 0, an observation due to B. White [37].

In fact, their reasoning can be considerably simplified by combining Theorem 4 with Moser's weak Harnack inequality for parabolic equations. First note:

**Theorem 6.** *For any  $T > 0$  there exists  $\varepsilon_1 > 0$  depending only on  $T, M$  and  $N$  such that any smooth solution  $u : M \times [0, T] \rightarrow N$  of (1.4), (1.5) with  $E(u_0) < \varepsilon_1$  can be extended to a global, smooth solution  $u : M \times [0, \infty] \rightarrow N$ , converging, as  $t \rightarrow \infty$  suitably, to a constant harmonic map.*

*Proof.* Let  $R_0^2 = \inf\{r_M^2, T\}$ . For  $0 < R_0^2 \leq t_0 \leq T$ ,  $x_0 \in M$  we can estimate with constants  $C$  depending on  $M, N$ , and  $T$  only

$$\Phi(R_0; z_0) \leq CR_0^{2-m}E(u(t_0 - R_0^2)) \leq CE(u_0) < \varepsilon_0, \quad (3.3)$$

if  $\varepsilon_1 > 0$  is sufficiently small. Hence by Theorem 4 we have

$$|\nabla u(x, t)| \leq C \text{ uniformly for } t \geq R_0^2, x \in M,$$

and  $u$  can be extended as a smooth solution of (1.4), (1.5) on  $M \times [0, \infty]$ .

By (1.8) and Moser's [28] supremum estimate for weak sub-solutions of linear parabolic equations, moreover we obtain

$$|\nabla u(x, t)|^2 \leq C E(u(t - R_0^2)) \leq CE(u_0) \quad \text{for } t \geq 2R_0^2, x \in M. \quad (3.4)$$

From this uniform estimate, asymptotic convergence follows as in Eells-Sampson [13] or Jost [21]. Finally, if  $\varepsilon_1 > 0$  is sufficiently small, by (3.4) the image of any map  $u(t)$ ,  $t \geq 2R_0^2$ , and hence also of the limiting harmonic map  $u_\infty$  is contained in a strictly convex geodesic ball on  $N$ . It follows that  $u_\infty \equiv \text{const.}$ ; see Jäger-Kaul [19].  $\square$

By Theorem 6, of course, for homotopically non-trivial initial data  $u_0$  with  $E(u_0) < \varepsilon_1(T)$  the flow (1.4), (1.5) must blow up before time  $T$ . In fact, blow-up time approaches 0 as the initial energy decreases to 0.

Finally, we remark that in dimensions  $m \geq 3$  singularities – as in the case  $m = 2$  – seem to be related to harmonic spheres or to self-similar solutions  $u(x, t) = w\left(\frac{x-x_0}{\sqrt{t_0-t}}\right)$  of (1.4); see Struwe [36; Theorem 8.1]. (The work of Coron-Ghidaglia strongly suggests that solutions of the latter kind in dimensions  $m \geq 3$  actually exist.)

The approach of [8] in general cannot be extended to initial maps belonging to  $H^{1,2}(M; N)$ , only. A different approach via time-discrete minimization was proposed by Horihata-Kikuchi [18]. Based on their ideas, Bethuel et al. [1] recently established global existence of distribution solutions to (1.4) for finite energy maps into spheres.

Further directions of present research into (1.4) include the study of (1.4) on complete, non-compact manifolds; see Li-Tam [25] for some recent work in this regard.

A subject related to (1.4) is the Cauchy problem for harmonic maps into Minkowski space. See Shatah [31], Sideris [32].

## References

1. Bethuel, F., Coron, J.M., Ghidaglia, J.M., Soyeur, A.: Heat flows and relaxed energies for harmonic maps. Preprint 1990
2. Brezis, H., Coron, J.M., Lieb, E.: Harmonic maps with defects. *Comm. Math. Phys.* **107** (1986) 649–705
3. Chang, K.-C.: Heat flow and boundary value problem for harmonic maps. Preprint 1988
4. Chang, K.-C., Ding, W.-Y.: A result on global existence for heat flows of harmonic maps from  $D^2$  into  $S^2$ . Preprint 1989
5. Chen, Y.: Weak solutions to the evolution problems of harmonic maps. *Math. Z.* **201** (1989) 69–74
6. Chen, Y., Ding, W.-Y.: Blow-up and global existence for heat flows of harmonic maps. Preprint
7. Chen, Y., Musina, R.: Harmonic mappings into manifolds with boundary. Preprint, Trieste 1989
8. Chen, Y., Struwe, M.: Existence and partial regularity results for the heat flow for harmonic maps. *Math. Z.* **201** (1989) 83–103
9. Cheng, X.: Estimate of singular set of the evolution problem for harmonic maps. Preprint, Rice Univ. 1990
10. Coron, J.-M.: Nonuniqueness for the heat flow of harmonic maps. Preprint 1988
11. Coron, J.-M., Ghidaglia, J.-M.: Explosion en temps fini pour le flot des applications harmoniques. Preprint 1988
12. Eells, J., Lemaire, L.: A report on harmonic maps. *Bull. London Math. Soc.* **10** (1978) 1–68
13. Eells, J., Sampson, J.H.: Harmonic mappings of Riemannian manifolds. *Amer. J. Math.* **86** (1964) 109–160
14. Eells, J., Wood, J.C.: Restrictions on harmonic maps of surfaces. *Topology* **15** (1976) 263–266
15. Giaquinta, M., Giusti, E.: On the regularity of the minima of variational integrals. *Acta Math.* **148** (1982) 31–46
16. Grayson, M., Hamilton, R.S.: The formation of singularities in the harmonic map heat flow. Preprint
17. Hamilton, R.S.: Harmonic maps of manifolds with boundary. (Lecture Notes in Mathematics, vol. 471.) Springer, Berlin Heidelberg New York 1975
18. Horihata, H., Kikuchi, N.: A construction of solutions satisfying a Caccioppoli inequality for nonlinear parabolic equations associated to a variational functional of harmonic type. Preprint
19. Jäger, W., Kaul, H.: Uniqueness and stability of harmonic maps, and their Jacobi fields. *Manuscr. math.* **28** (1979) 269–291
20. Jäger, W., Kaul, H.: Rotationally symmetric harmonic maps from a ball into a sphere and the regularity problem for weak solutions of elliptic systems. *J. Reine Angew. Math.* **343** (1983) 146–161
21. Jost, J.: Ein Existenzbeweis für harmonische Abbildungen, die ein Dirichletproblem lösen, mittels der Methode des Wärmeflusses. *Manuscr. math.* **34** (1981) 17–25
22. Keller, J., Rubinstein, J., Sternberg, P.: Reaction – diffusion processes and evolution to harmonic maps. Preprint
23. Lemaire, L.: Applications harmoniques de surfaces riemannniennes. *J. Diff. Geom.* **13** (1978) 51–78
24. Li, M.: Harmonic map heat flow with free boundary. Preprint, Trieste 1990
25. Li, P., Tam, L.-F.: The heat equation and harmonic maps of complete manifolds. Preprint 1989

26. Lin, F.H.: Une remarque sur l'application  $x/|x|$ . C.R. Acad. Sc. Paris **305** (1987) 529–531
27. Mitteau, J.C.: Sur les applications harmoniques. J. Diff. Geom. **9** (1974) 41–54
28. Moser, J.: A Harnack inequality for parabolic differential equations. Comm. Pure Appl. Math. **17** (1964) 101–134
29. Sacks, J., Uhlenbeck, K.: The existence of minimal immersions of 2-spheres. Ann. Math. **113** (1981) 1–24
30. Schoen, R.S., Uhlenbeck, K.: A regularity theory for harmonic maps. J. Diff Geom. **17** (1982) 307–335, and **18** (1983) 329
31. Shatah, J.: Weak solutions and the development of singularities of the  $SU(2)$   $\sigma$ -model. Comm. Pure Appl. Math. **41** (1988) 459–469
32. Sideris, T.: Global existence of harmonic maps in Minkowski space. Comm. Pure Appl. Math. **42** (1989) 1–13
33. Struwe, M.: On the evolution of harmonic maps of Riemannian surfaces. Comment. Math. Helv. **60** (1985) 558–581
34. Struwe, M.: The existence of surfaces of constant mean curvature with free boundaries. Acta Math. **160** (1988) 19–64
35. Struwe, M.: Heat flow methods for harmonic maps of surfaces and applications to free boundary problems. (Lecture Notes in Mathematics, vol. 1324.) Springer, Berlin Heidelberg New York 1988, pp. 293–319
36. Struwe, M.: On the evolution of harmonic maps in higher dimensions. J. Diff. Geom. **28** (1988) 485–502
37. White, B.: Infima of energy functionals in homotopy classes of mappings. J. Diff. Geom. **23** (1986) 127–142



# Integrable Systems in Gauge Theory, Kähler Geometry and Super KP Hierarchy – Symmetries and Algebraic Point of View

*Kanehisa Takasaki*

Research Institute for Mathematical Sciences, Kyoto University  
Sakyo-ku, Kyoto 606, Japan

## 1. Introduction

The following nonlinear systems all provide valuable material to search for new “nonlinear integrable systems”.

- self-duality equation in Yang-Mills theory
- self-duality equation in Kähler geometry
- super Kadomtsev-Petviashvili (KP) hierarchy.

From these equations, one will be able to imagine several types of extensions of so called “soliton equations” such as the celebrated Korteweg-de Vries (KdV) equation etc. The first two cases are in a sense “higher dimensional” (or “multi-dimensional”) nonlinear integrable systems; the last case will be interesting as an extension of M. Sato’s work on the KP hierarchy [SS] and background ideas [S] referred to under the key words “algebraic analysis.”

This lecture is a summary of my recent work on these equations, in particular, the self-duality equations, with focus on their symmetry properties. It is nowadays widely recognized that symmetries of soliton equations can be described by representation theory of Kac-Moody Lie algebras [DJKM]. A similar observation to the self-duality equation of Yang-Mills theory has been known for years [UN, CGW, D, T1]. The case of the self-duality equation in Kähler geometry seems to have remained less obscure [BP]. Recently I obtained an explicit description of infinitesimal symmetries, which exhibits a Poisson algebra structure [T2]. Very recently, inspired by a work of Leznov et al. [LMS], I noticed that these infinitesimal symmetries can be “exponentiated” by a simple method [T3]. This leads to a kind of “perturbative” construction of a class of general (local) solutions. To stress underlying symplectic structures, I will illustrate these results for a  $4N$ -dimensional generalization of the self-duality equations rather than in the original form.

The basic standpoint of my work largely relies on the philosophy of “algebraic analysis,” which understands differential equations as a differential algebra, i.e., a set of abstract symbols and differential-algebraic relations among them. This language has turned out to be particularly useful [T4] in the case of the super KP hierarchy of Manin and Radul [MR] as well as the original KP hierarchy. For the treatment of the self-duality equations, we shall not specify such a differential-algebraic interpretation; however, its spirit is included therein.

## 2. Generalized Self-Duality Equations

### 2.1 The Case of Yang-Mills Theory

We consider a  $4N$ -dimensional space-time with coordinates

$$(x, p) = (x^1, \dots, x^{2N}, p^1, \dots, p^{2N}) \quad (1)$$

and a generalized self-duality equation of Yang-Mills theory on this space-time. This equation, as in the four dimensional case, has two equivalent expressions [C]. As we shall see later on, these two expressions have analogues in Kähler geometry. The first expression is given by the equations

$$\frac{\partial^2 K}{\partial x^i \partial p^j} - \frac{\partial^2 K}{\partial x^j \partial p^i} + \left[ \frac{\partial K}{\partial x^i}, \frac{\partial K}{\partial x^j} \right]^\circ = 0, \quad (2)$$

where the unknown function  $K = K(x, p)$  takes values in the Lie algebra  $\text{Lie}G$  of the structure group  $G$ . The second one is given by

$$\frac{\partial}{\partial x^i} \left( \frac{\partial J}{\partial p^j} J^{-1} \right) - \frac{\partial}{\partial x^j} \left( \frac{\partial J}{\partial p^i} J^{-1} \right) = 0, \quad (3)$$

where the unknown function  $J = J(x, p)$  now takes values in  $G$ .

As well known, these equations are the integrability condition (in the sense of Frobenius) of the linear system

$$\left( \frac{\partial}{\partial p^i} - \lambda \frac{\partial}{\partial x^i} + A_i \right) \Psi(\lambda) = 0. \quad (4)$$

The gauge potentials  $A_i$  are combined with the previous unknown functions  $J$  and  $K$  as:

$$A_i = -\frac{\partial K}{\partial x^i} = -\frac{\partial J}{\partial p^i} J^{-1}. \quad (5)$$

We consider, in particular, a special pair of solutions

$$\begin{aligned} \Psi(\lambda) &= W(\lambda), \quad W(\lambda) = 1 + \sum_{n \leq -1} W_n \lambda^n, \\ \Psi(\lambda) &= V(\lambda), \quad V(\lambda) = \sum_{n \geq 0} V_n \lambda^n \end{aligned} \quad (6)$$

connected with  $J$  and  $K$  by the relation

$$K = -W_{-1}, \quad J = V_0. \quad (7)$$

The linear system, with these expressions inserted, gives rise to a *nonlinear* system with the new unknown functions  $W_n$  and  $V_n$ . Symmetries are to be constructed for this nonlinear system rather than the original equation.

## 2.2 The Case of Kähler Geometry

We now turn to Kähler geometry. Our notational conventions are as follows. Let  $i, j, \dots$  be symplectic indices with values in integers  $1, \dots, 2N$ .  $\epsilon^{ij}$  and  $\epsilon_{ij}$  denote the standard symplectic  $\epsilon$ -symbols normalized as  $\epsilon_{2i-1, 2i} = -\epsilon_{2i, 2i-1} = 1$  and  $\epsilon^{2i-1, 2i} = -\epsilon^{2i, 2i-1} = 1$ . The Einstein summation convention is understood only for symplectic indices. Symplectic indices are raised and lowered as  $\xi_i = \epsilon_{ij} \xi^j$  and  $\eta^j = \eta_i \epsilon_{ij}$ .

A  $4N$ -dimensional generalization of the self-duality equation in Kähler geometry is provided by hyper-Kähler geometry. As pointed out (or re-discovered) by Plebanski [P] in the four dimensional (self-dual) case, hyper-Kähler geometry (also called “ $\mathcal{H}$ -space”) has two equivalent local pictures based upon the first and second “heavenly equations.” The “second” picture consists of a  $4N$ -dimensional coordinate system  $(x, p) = (x^1, \dots, x^{2N}, p^1, \dots, p^{2N})$ , a scalar unknown function  $\Theta = \Theta(x, p)$ , and the “second heavenly equation”

$$\frac{\partial^2 \Theta}{\partial x^i \partial p^j} - \frac{\partial^2 \Theta}{\partial x^j \partial p^i} + \left\{ \frac{\partial \Theta}{\partial x^i}, \frac{\partial \Theta}{\partial x^j} \right\}_{(x)} = 0, \quad (8)$$

where  $\{ \ , \ \}_{(x)}$  stands for the Poisson bracket in  $x$ ,

$$\{F, G\}_{(x)} = \epsilon^{ij} \frac{\partial F}{\partial x^i} \frac{\partial G}{\partial x^j}. \quad (9)$$

In the “first” picture, one has a  $4N$ -dimensional coordinate system  $(p, \hat{p}) = (p^1, \dots, p^{2N}, \hat{p}^1, \dots, \hat{p}^{2N})$ , a scalar unknown function  $\Omega = \Omega(p, \hat{p})$ , and the “first heavenly equation”

$$\left\{ \frac{\partial \Omega}{\partial p^i}, \frac{\partial \Omega}{\partial \hat{p}^j} \right\}_{(\hat{p})} = \epsilon_{ij}, \quad (10)$$

where we use another Poisson bracket,

$$\{F, G\}_{(\hat{p})} = \epsilon^{ij} \frac{\partial F}{\partial \hat{p}^i} \frac{\partial G}{\partial \hat{p}^j}. \quad (11)$$

Geometrically,  $\Omega$  represents a Kähler potential, and  $p^i$  and  $\hat{p}_i$  correspond to complex coordinates and their complex conjugates. In the following, however, we understand  $(p, \hat{p})$  or  $(x, p)$  as  $4N$  independent complex variables and consider formal aspects of the above differential equations.

The role of  $W(\lambda)$  and  $V(\lambda)$  is now to be played by two sets of functions (or formal Laurent series)

$$\begin{aligned} u^i(\lambda) &= \sum_{n \leq -1} u_n^i \lambda^n + x^i + p^i \lambda \quad (1 \leq i \leq 2N), \\ \hat{u}^i(\lambda) &= \hat{p}^i + \sum_{n \geq 1} \hat{u}_n^i \lambda^n, \quad (1 \leq i \leq 2N) \end{aligned} \quad (12)$$

subject to the exterior differential equations

$$\epsilon_{ij} du^i(\lambda) \wedge d\bar{u}^j(\lambda) = \epsilon_{ij} d\hat{u}^i(\lambda) \wedge d\bar{\hat{u}}^j(\lambda), \quad (13)$$

and

$$\begin{aligned} d\Theta &= \in_{ij} u^i_{-2} dp^j + \in_{ij} u^i_{-1} dx^i, \\ d\Omega &= - \in_{ij} u^i_0 dp^j + \in_{ij} \hat{u}^i_1 d\hat{p}^j. \end{aligned} \quad (14)$$

Here  $u_n^i$  and  $\hat{u}_n^i$  are understood as unknown functions of  $(x, p)$  (in the second heavenly picture) or of  $(p, \hat{p})$  (in the first heavenly picture);  $\lambda$  is considered a constant under the total differential  $d$ , i.e.,  $d\lambda = 0$ . Symmetries are to be constructed for this nonlinear system.

### 3. Infinitesimal Symmetries

#### 3.1 The Case of Yang-Mills Theory [T1]

For the  $(W(\lambda), V(\lambda))$ -system, a one-parameter family of transformations

$$(W(\lambda), V(\lambda)) \mapsto (W(\varepsilon, \lambda), V(\varepsilon, \lambda)) \quad (15)$$

of solutions is defined by the Riemann-Hilbert factorization

$$W(\varepsilon, \lambda) e^{-\varepsilon X(\lambda)} W(\lambda)^{-1} = V(\varepsilon, \lambda) e^{-\varepsilon Y(\lambda)} V(\lambda)^{-1}. \quad (16)$$

Here  $X(\lambda) = X(\lambda, x, p)$  and  $Y(\lambda) = Y(\lambda, x, p)$ , the data of transformations, are LieG-valued functions of  $4N + 1$  variables of the form

$$\begin{aligned} X(\lambda) &= \mathbf{X}(\lambda, x^1 + p^1 \lambda, \dots, x^{2N} + p^{2N} \lambda), \\ Y(\lambda) &= \mathbf{Y}(\lambda, x^1 + p^1 \lambda, \dots, x^{2N} + p^{2N} \lambda), \end{aligned} \quad (17)$$

where  $\mathbf{X}$  and  $\mathbf{Y}$  are arbitrary LieG-valued functions of  $2N + 1$  variables with Laurent expansion

$$\mathbf{X}(\lambda, u) = \sum_{n=-\infty}^{\infty} \mathbf{X}_n(u) \lambda^n, \quad \mathbf{Y}(\lambda, u) = \sum_{n=-\infty}^{\infty} \mathbf{Y}_n(u) \lambda^n. \quad (18)$$

[In fact, some restriction on these data is required for the Riemann-Hilbert factorization to work well; a prescription is to put upper and lower bounds to the range of  $n$  as  $-\infty < n \leq n_X$  for  $\mathbf{X}(\lambda)$  and  $-n_Y \leq n < \infty$  for  $\mathbf{Y}(\lambda)$ . A similar remark also applies to the hyper-Kähler case. This is a somewhat technical issue.] The associated infinitesimal transformations

$$\begin{aligned} \delta_{X,Y} W(\lambda) &= \left. \frac{\partial W(\varepsilon, \lambda)}{\partial \varepsilon} \right|_{\varepsilon=0}, \\ \delta_{X,Y} V(\lambda) &= \left. \frac{\partial V(\varepsilon, \lambda)}{\partial \varepsilon} \right|_{\varepsilon=0} \end{aligned} \quad (19)$$

have the following structure.

**Proposition 1.** *The infinitesimal symmetries act on  $W(\lambda)$  and  $V(\lambda)$  as follows.*

$$\begin{aligned} \delta_{X,Y} W(\lambda) \cdot W(\lambda)^{-1} &= (W(\lambda) X(\lambda) W(\lambda)^{-1} - V(\lambda) Y(\lambda) V(\lambda)^{-1})_{\leq -1}, \\ \delta_{X,Y} V(\lambda) \cdot V(\lambda)^{-1} &= (V(\lambda) Y(\lambda) V(\lambda)^{-1} - W(\lambda) X(\lambda) W(\lambda)^{-1})_{\geq 0}, \end{aligned} \quad (20)$$

where  $(\ )_{\geq 0}$  and  $(\ )_{\leq -1}$  are linear maps on the space of Laurent series of  $\lambda$  defined by

$$\begin{aligned} (\ )_{\geq 0} : \lambda^n &\mapsto \theta(n \geq 0) \lambda^n, \\ (\ )_{\leq -1} : \lambda^n &\mapsto \theta(n \leq -1) \lambda^n. \end{aligned} \quad (21)$$

Further, these infinitesimal symmetries obey the commutation relations

$$[\delta_{X_1, Y_1}, \delta_{X_2, Y_2}] = \delta_{[X_1, X_2], [Y_1, Y_2]}. \quad (22)$$

Thus, in particular, the infinitesimal symmetries give rise to a nonlinear realization of a direct sum of two loop algebras (with extra  $2N$  variables  $u^1, \dots, u^{2N}$ ). The associated infinitesimal transformations of  $J = V_0$  and  $K = -W_{-1}$  can be readily derived from the above result.

### 3.2 The Case of Kähler Geometry [T2]

The case of  $(u(\lambda), \hat{u}(\lambda))$ -system requires a more involved factorization, i.e., a factorization with respect to composition of maps. Let us consider this issue within the  $(x, p)$ -coordinate system. [A fully parallel treatment is possible with the  $(p, \hat{p})$ -coordinate system.]  $u(\lambda)$  and  $\hat{u}(\lambda)$  are now interpreted as maps

$$\begin{aligned} u(\lambda) : x &\mapsto u(\lambda, x, p), \\ \hat{u}(\lambda) : x &\mapsto \hat{u}(\lambda, x, p) \end{aligned} \quad . \quad (23)$$

from the  $x$ -space into the  $u$ -space or  $\hat{u}$ -space respectively. A one-parameter family of transformations

$$(u(\lambda), \hat{u}(\lambda)) \mapsto (u(\varepsilon, \lambda), \hat{u}(\varepsilon, \lambda)) \quad (24)$$

of solutions can be defined by the Riemann-Hilbert factorization

$$u(\varepsilon, \lambda)^{-1} \circ e^{-\varepsilon \xi_F(\lambda)} \circ u(\lambda) = \hat{u}(\varepsilon, \lambda)^{-1} \circ e^{-\varepsilon \xi_{\hat{F}}(\lambda)} \circ \hat{u}(\lambda), \quad (25)$$

where  $\xi_F(\lambda)$  and  $\xi_{\hat{F}}(\lambda)$  are Hamiltonian vector fields of the form

$$\begin{aligned} \xi_F(\lambda) &= \epsilon^{ij} \frac{\partial F(\lambda)}{\partial u^i} \frac{\partial}{\partial u^j}, \\ \xi_{\hat{F}}(\lambda) &= \epsilon^{ij} \frac{\partial \hat{F}(\lambda)}{\partial \hat{u}^i} \frac{\partial}{\partial \hat{u}^j}. \end{aligned} \quad (26)$$

The generating functions  $F(\lambda) = F(\lambda, u)$  and  $\hat{F}(\lambda) = \hat{F}(\lambda, \hat{u})$  are arbitrary functions of  $2N + 1$  variables with Laurent expansion

$$F(\lambda) = \sum_{n=-\infty}^{\infty} F_n(u) \lambda^n, \quad \hat{F}(\lambda) = \sum_{n=-\infty}^{\infty} \hat{F}_n(\hat{u}) \lambda^n. \quad (27)$$

The infinitesimal transformations

$$\begin{aligned}\delta_{F,\hat{F}} u^i(\lambda) &= \left. \frac{\partial u^i(\varepsilon, \lambda)}{\partial \varepsilon} \right|_{\varepsilon=0}, \\ \delta_{F,\hat{F}} \hat{u}^i(\lambda) &= \left. \frac{\partial \hat{u}^i(\varepsilon, \lambda)}{\partial \varepsilon} \right|_{\varepsilon=0}\end{aligned}\quad (28)$$

have the following structure.

**Proposition 2.** *The infinitesimal symmetries act on  $u(\lambda)$  and  $\hat{u}(\lambda)$  as:*

$$\begin{aligned}\delta_{F,\hat{F}} u^i(\lambda) &= \left\{ \left[ F(\lambda, u(\lambda)) - \hat{F}(\lambda, \hat{u}(\lambda)) \right]_{\leq -1}, u^i(\lambda) \right\}_{(x)}, \\ \delta_{F,\hat{F}} \hat{u}^i(\lambda) &= \left\{ \left[ \hat{F}(\lambda, \hat{u}(\lambda)) - F(\lambda, u(\lambda)) \right]_{\geq 0}, \hat{u}^i(\lambda) \right\}_{(x)}.\end{aligned}\quad (29)$$

Further, the infinitesimal symmetries obey the commutation relations

$$\left[ \delta_{F_1, \hat{F}_1}, \delta_{F_2, \hat{F}_2} \right] = \delta_{\{F_1, F_2\}_{(0)}, \{\hat{F}_1, \hat{F}_2\}_{(0)}}. \quad (30)$$

Thus the infinitesimal symmetries give a nonlinear realization of a direct sum of two Poisson (loop) algebras.

Remarkably, the above infinitesimal symmetries can be extended to  $\Theta$  and  $\Omega$  without modifying the Poisson algebra structure.

**Proposition 3.** *The infinitesimal symmetries can be consistently extended to  $\Theta$  and  $\Omega$  by the following rule.*

$$\begin{aligned}\delta_{F,\hat{F}} \Theta &= \text{res}_{\lambda=\infty} F(\lambda, u(\lambda)) + \text{res}_{\lambda=0} \hat{F}(\lambda, \hat{u}(\lambda)), \\ \delta_{F,\hat{F}} \Omega &= - \text{res}_{\lambda=\infty} \lambda^{-2} F(\lambda, u(\lambda)) - \text{res}_{\lambda=0} \lambda^{-2} \hat{F}(\lambda, \hat{u}(\lambda)),\end{aligned}\quad (31)$$

where the residues are normalized as

$$\text{res}_{\lambda=\infty} \lambda^n = -\delta_{n,-1}, \quad \text{res}_{\lambda=0} \lambda^n = \delta_{n,-1}. \quad (32)$$

These extended infinitesimal symmetries obey the same commutation relations as in Proposition 2.

#### 4. Perturbative Method [T3]

The infinitesimal symmetries, as we have seen, have a considerably simple and beautiful structure. The Riemann-Hilbert factorization problems in general are hard to solve explicitly. For the case of Yang-Mills fields, several solution methods are developed; for the hyper-Kähler case, only existence theorems are known (except for a few very special families of solutions). The method presented here, so to speak, “exponentiate” the infinitesimal symmetries by expanding everything in powers of  $\varepsilon$ . As Leznov et al. [LMS] pointed out, the parameter  $\varepsilon$  plays the role of “coupling constants” in field theory; therefore we call the following method “perturbative.”

#### 4.1 The Case of Yang-Mills Theory

Let us consider the previous Riemann-Hilbert factorization in case where

$$W(\lambda) = V(\lambda) = 1 \text{ (trivial solution)}, \quad Y(\lambda) = 0. \quad (33)$$

Let us define

$$\mathcal{X}(\varepsilon, \lambda) = W(\varepsilon, \lambda)X(\lambda)W(\varepsilon, \lambda)^{-1} \quad (34)$$

Since  $\partial/\partial\varepsilon$  corresponds to the action of  $\delta_{X,0}$  on  $(W(\varepsilon, \lambda), V(\varepsilon, \lambda))$ , one can readily find a closed differential equation satisfied by  $\mathcal{X}(\varepsilon, \lambda)$  with respect to  $\varepsilon$ .

**Proposition 4.**  $\mathcal{X}(\varepsilon, \lambda)$  satisfies the differential equation

$$\frac{\partial \mathcal{X}(\varepsilon, \lambda)}{\partial \varepsilon} = \left[ (\mathcal{X}(\varepsilon, \lambda))_{\leq -1}, \mathcal{X}(\varepsilon, \lambda) \right] \quad (35)$$

and the initial condition

$$\mathcal{X}(\varepsilon = 0, \lambda) = \mathbf{X}(\lambda, x + p\lambda). \quad (36)$$

Further,

**Proposition 5.**  $K(\varepsilon) = -W_{-1}(\varepsilon)$  and  $J(\varepsilon) = V_0(\varepsilon)$  obey the differential equations

$$\begin{aligned} \frac{\partial K(\varepsilon)}{\partial \varepsilon} &= \operatorname{res}_{\lambda=\infty} \mathcal{X}(\varepsilon, \lambda), \\ \frac{\partial J(\varepsilon)}{\partial \varepsilon} J(\varepsilon)^{-1} &= \operatorname{res}_{\lambda=\infty} \lambda^{-1} \mathcal{X}(\varepsilon, \lambda). \end{aligned} \quad (37)$$

Substitution of the Taylor expansion (“perturbation series”)

$$\mathcal{X}(\varepsilon, \lambda) = \sum_{k=0}^{\infty} \mathcal{X}^{(k)}(\lambda) \varepsilon^k / k! \quad (38)$$

into the above equation yields a set of recursive relations

$$\begin{aligned} \mathcal{X}^{(0)}(\lambda) &= X(\lambda) = \mathbf{X}(\lambda, x + p\lambda), \\ \mathcal{X}^{(k+1)} &= \sum_{\ell=0}^k \binom{k}{\ell} \left[ (\mathcal{X}^{(k-\ell)}(\lambda))_{\leq -1}, \mathcal{X}^{(\ell)}(\lambda) \right]. \end{aligned} \quad (39)$$

The unknown functions  $K(\varepsilon)$  and  $J(\varepsilon)$  of the generalized self-duality equations, too, can be determined by expansion into powers of  $\varepsilon$ .

In the original formulation of Leznov et al. [LMS], the projection  $(\ )_{\leq -1}$  is represented by an integral operator; they exploit its algebraic properties to check, by brute force, the validity of their formula.

### 3.2 The Case of Kähler Geometry

We now start from the Riemann-Hilbert factorization with

$$u^i(\lambda) = \hat{u}^i(\lambda) = x^i + p^i\lambda \text{ (trivial solution)}, \quad \hat{F}(\lambda) = 0, \quad (40)$$

and derive differential equations satisfied by

$$\mathcal{F}(\varepsilon, \lambda) = F(\lambda, u(\varepsilon, \lambda)) \quad (41)$$

and  $\Theta(\varepsilon)$  with respect to  $\varepsilon$ .

**Proposition 6.**  $\mathcal{F}(\varepsilon, \lambda)$  satisfies the differential equation

$$\frac{\partial \mathcal{F}(\varepsilon, \lambda)}{\partial \varepsilon} = \{[\mathcal{F}(\varepsilon, \lambda)]_{\leq -1}, \mathcal{F}(\varepsilon, \lambda)\}_{(x)} \quad (42)$$

and the initial condition

$$\mathcal{F}(\varepsilon = 0, \lambda) = F(\lambda, x + p\lambda). \quad (43)$$

**Proposition 7.** One can obtain  $\Theta(\varepsilon)$  by solving the differential equation

$$\frac{\partial \Theta(\varepsilon)}{\partial \varepsilon} = \underset{\lambda \rightarrow \infty}{\text{res}} \mathcal{F}(\varepsilon, \lambda) \quad (44)$$

under the initial condition

$$\Theta(\varepsilon = 0) = 0. \quad (45)$$

These equations can be solved by the same ‘‘perturbative method’’ as illustrated in the case of Yang-Mills fields.

The above construction is not suited for the first heavenly picture based upon  $(p, \hat{p}, \Omega)$ . To give a similar construction for the first heavenly picture, we just have to restart from the situation where

$$u^i(\lambda) = \hat{u}^i(\lambda) = \hat{p}^i + p^i\lambda, \quad F(\lambda) = 0, \quad (46)$$

and consider equations satisfied by  $\Omega(\varepsilon)$  and

$$\hat{\mathcal{F}}(\varepsilon, \lambda) = \hat{F}(\lambda, \hat{u}(\varepsilon, \lambda)). \quad (47)$$

## 5. KP and Super KP Hierarchies [T4]

In the differential-algebraic approach mentioned in the introduction, a nonlinear system is represented by a commutative algebra  $\mathcal{A}$  with a set of derivations  $\partial_1, \partial_2, \dots$ . If one is not interested in a particular choice of such derivations, it is convenient to understand a differential algebra as a pair  $(\mathcal{A}, \Delta)$  of a commutative algebra and an  $\mathcal{A}$ -module  $\Delta$  of derivations in  $\mathcal{A}$ . Infinitesimal symmetries are then, by definition, derivations  $\delta : \mathcal{A} \rightarrow \mathcal{A}$  that satisfy the condition

$$[\delta, \partial] \in \Delta \quad (\forall \partial \in \Delta). \quad (48)$$

In most applications, the derivations  $\delta_1, \delta_2, \dots$ , are chosen to be commutative, and symmetries are characterized as extra derivations of  $\mathcal{A}$  that commute with these derivations. (The super KP hierarchy is somewhat distinct; not only  $\mathcal{A}$  being a supercommutative algebra, the set of derivations are neither commutative nor supercommutative. One can however see that its basic structure is almost parallel to the case of the KP hierarchy.)

The generalized self-duality equations, too, can be formulated as such an abstract differential algebra. Its algebraic part  $\mathcal{A}$  should be a commutative algebra (over a suitable differential subalgebra that specifies in which domain to seek for solutions) generated by the Laurent coefficients of  $W(\lambda)$  and  $V(\lambda)$ , or of  $u(\lambda)$  and  $\hat{u}(\lambda)$ . In the latter case, one may also add  $\Theta$  or  $\Omega$ . This certainly provides an unambiguous framework for the notion of infinitesimal symmetries; however, one will gain nothing practically new from this reinterpretation.

The situation is considerably different for the case of the KP and super KP hierarchies. For these equations, the differential-algebraic language seems to have a substantial meaning. Of particular importance is a  $\mathcal{D}$ -module structure hidden in the formulation of the KP and super KP hierarchy. (This observation for the case of the KP hierarchy is due to Sato, who stresses the relevance of the notion of  $\mathcal{D}$ -modules even in more general perspectives [S].) With the aid of this  $\mathcal{D}$ -module structure, one can find a new set of generators  $w_{ij}$  ( $i \geq 0, j \leq -1$ ) in  $\mathcal{A}$ . This is the most direct way to see a connection with the geometry of infinite dimensional (super) Grassmannian manifolds;  $w_{ij}$ 's can be identified with affine coordinates on an open subset therein. This also leads to: an explicit description of infinitesimal symmetries  $\delta_A$  parametrized by elements  $A$  of an infinite matrix Lie algebra  $gl(\infty)$  (for the KP hierarchy) or of its super-version  $gl(\infty|\infty)$  (for the super KP hierarchy), a differential-algebraic characterization of the  $\tau$  function, its symmetry contents related to central extensions of  $gl(\infty)$  and  $gl(\infty|\infty)$ , etc.

In fact,  $\Theta$  and  $\Omega$  may be in a sense considered an analogue of the  $\tau$  function. This analogy becomes quite reasonable if we consider a hierarchy of the generalized self-duality equations discussed here. Their representation-theoretic and geometric properties are however considerably different from the  $\tau$  function.

## References

- [BP] Boyer, C.P., Plebanski, J.F.: An infinite hierarchy of conservation laws and nonlinear superposition principles for self-dual Einstein spaces. *J. Math. Phys.* **26** (1985) 229–234
- [C] Chau, L.-L.: Chiral fields, self-dual Yang-Mills fields as integrable systems, and the role of the Kac-Moody algebra. In: Nonlinear phenomena, K.B. Wolf (ed.). (Lecture Notes in Physics, vol. 189.) Springer, Berlin Heidelberg New York 1983
- [CGW] Chau, L.-L., Ge, M.-L., Wu, Y.-S.: Kac-Moody algebra in the self-dual Yang-Mills equation. *Phys. Rev. D* **25** (1982) 1086–1094
- [D] Dolan, L.: A new symmetry group of real self-dual Yang-Mills theory. *Phys. Lett.* **113B** (1982) 387–390
- [DJKM] Date, E., Jimbo, M., Kashiwara, M., Miwa, T.: Transformation theory for soliton equations III-VI. *J. Phys. Soc. Japan* **50** (1982) 3806–3812, 3813–3818; *Physica* **4D** (1982) 343–365. *Publ. RIMS, Kyoto Univ.* **18** (1982) 1077–1110

- [LMS] Leznov, A.N., Mukhtarov, M.A.: Deformation of algebras and solutions of self-duality equation. *J. Math. Phys.* **28** (1987) 2574–2578. Leznov, A.N., Saveliev, V.M.: Exactly and completely integrable nonlinear dynamical systems. *Acta Appl. Math.* **16** (1989) 1–74
- [MR] Manin, Yu.I., Radul, A.O.: A supersymmetric extension of the Kadomtsev-Petviashvili hierarchy. *Commun. Math. Phys.* **98** (1985) 65–77
- [P] Plebanski, J.F.: Some solutions of complex Einstein equations. *J. Math. Phys.* **16** (1975) 2395–2402
- [SS] Sato, M., Sato, Y.: Soliton equations as dynamical systems in an infinite dimensional Grassmann manifold. In: *Nonlinear Partial Differential Equations in Applied Sciences*, P.D. Lax, H. Fujita, G. Strang (eds.). North-Holland, Amsterdam, and Kinokuniya, Tokyo, 1982
- [S] Sato, M.:  $\mathcal{D}$ -modules and nonlinear systems. In: *Integrable Systems in Quantum Theory and Statistical Mechanics*, M. Jimbo, T. Miwa, A. Tsuchiya (eds.). *Adv. Stud. Pure Math.*, vol. 19, pp. 417–434. Kinokuniya, Tokyo 1989
- [T1] Takasaki, K.: A new approach to the self-dual Yang-Mills equations. *Commun. Math. Phys.* **94** (1984) 35–59. Hierarchy structure in integrable systems of gauge fields and underlying Lie algebras. *Commun. Math. Phys.* **127** (1990) 225–238
- [T2] Takasaki, K.: An infinite number of hidden variables in hyper-Kähler metrics. *J. Math. Phys.* **30** (1989) 1515–1521. Symmetries of hyper-Kähler (or Poisson gauge field) hierarchy. *J. Math. Phys.* **31** (1990) 1877–1888
- [T3] Takasaki, K.: Perturbative approach to self-duality equations. In preparation
- [T4] Takasaki, K.: Symmetries of the super KP hierarchy. *Lett. Math. Phys.* **17** (1989) 351–357. Differential algebras and  $\mathcal{D}$ -modules in super Toda lattice hierarchy. *Lett. Math. Phys.* **19** (1990) 229–236
- [UN] Ueno, K., Nakamura, Y.: Transformation theory for anti-self-dual equations and the Riemann-Hilbert problem. *Phys. Lett.* **109B** (1982) 273–278

# H-Measures and Applications

*Luc Tartar*

Department of Mathematics, Carnegie Mellon University, Pittsburgh, PA 15213, USA

## Introduction

Historians interested in the evolution of Science will probably be very surprised when they will analyse all the strange fashions that have struck the scientific community in the second part of this century. They will certainly be aware of the political orientation of those who had launched many of these fake new ideas and wonder why mathematicians had not been more rational in their behaviour. Will they find another example of political censorship than the one that had suppressed two pages of the introduction of an article where I had described my scientific ideas [1]? Will everyone agree that it was indeed slanderous that I had thanked there two of my teachers, Laurent Schwartz and Jacques-Louis Lions and that this part should definitely be cut? Had the censors thought that they could suppress the mention of other political facts and avoid me describing them elsewhere [2]? Historians may wonder at the stupidity of these censors and ponder if they had even understood the meaning of the few lines that they had spared at the beginning: “Il y a une différence énorme entre l'étude des singularités d'équations aux dérivées partielles (linéaires ou non) et celle de leurs oscillations: c'est la différence entre la physique classique et la physique quantique”.

In this article I had described my point of view that the study of oscillating solutions of partial differential equations was the key mathematical question to investigate in order to shed some light on the strange rules invented by physicists for explaining natural phenomena. Every specialist of differential equations is aware of the distinction between finite and infinite dimensional effects and it is only the result of an intensive propaganda that so many have adopted a point of view about mechanics which was adequate in the eighteenth century when partial differential equations had not yet found their place and that continuum mechanics and electromagnetism were not even thought of. However, even if all the extensive knowledge about linear partial differential equations contained in the treatise of L. Hörmander [3] had been available at the beginning of the century, it would not have helped so much the physicists puzzled as they were by the spectroscopic measurements of light absorbed and emitted in some gases. One cannot blame then those who have invented the strange rules of quantum

---

<sup>1</sup> There is a huge difference between the study of singularities of partial differential equations (be them linear or not) and that of their oscillations: it is the difference between classical and quantum physics.

mechanics for their lack of knowledge of partial differential equations but one should blame those who have transformed these rules into dogma. As R. Penrose once wrote “Quantum theory, it may be said, has two things in its favour and only one against. First, it agrees with all the experiments. Second, it is a theory of astonishing and profound mathematical beauty. The only thing to be said against the theory is that it makes absolutely no sense”<sup>2</sup>. In order to give a rational explanation of the puzzling measurements made in these experiments one needs an increased knowledge of partial differential equations and there are a few new mathematical questions that should be understood for that purpose.

As light is involved we know that there will be some hyperbolic equations, the wave equation or Maxwell’s system or some even larger system, and we expect that the linearised system will only have the velocity of light as characteristic speed so that we may reasonably restrict our attention to semilinear systems. Even standard questions as the relation between the wave equation and geometrical optics needs to be thought again. Fifteen years ago it was already clear why the mathematical results now found in [3] were not adapted to the goal that I was pointing at and a first reason was that one cannot expect to understand the partial differential equations of continuum mechanics without accepting discontinuous coefficients; even if one was ready to make smoothness assumptions and stay away from interfaces one definitely had to avoid assuming the coefficients to be infinitely differentiable or analytic. There is however a more important drawback of the linear theory of propagation of singularities which became more apparent once I had obtained my personal version of propagation using the tool of H-measures [4] which I will describe in a moment. What the linear theory is really interested in is the propagation of regularity and this leads to a quite negative concept of a singularity which is not defined as a quantitative object; of course, measuring the propagation of  $H^s$  regularity instead of  $C^\infty$  regularity does not correct this defect in any way. The physically intuitive idea of a beam of light is then absolutely not described by the theory of “propagation of singularities” for partial differential equations and the inadequacy is hidden by the fact that the bicharacteristic rays which have appeared in the linear theory are precisely those which physicists had thought important in their formal computations. One should then criticise this approach of propagation of singularities for describing the properties of light as not making more sense than some physicists’ rules; a better test for a mathematical tool than making the bicharacteristic rays appear is to be able to measure what is transported along them and tell what happens along the bicharacteristic rays to important quantities for physics like energy and momentum.

As matter is also involved we face much more trouble because the question of what matter could be is at stake anyway and, even if the rules of quantum mechanics are indeed wrong, one cannot forget about the real defects of the classical concepts of light and matter. A probably good mathematical model to understand is the coupled Maxwell-Dirac system where matter is described by a complex four dimensional vector field and light is described by the electromagnetic field, the coupling involving quadratic terms with the famous Planck constant  $\hbar$  appearing as a coupling parameter between light and matter and not as this mysterious parameter that the dogma wants to attach to every hamiltonian.

---

<sup>2</sup> This is the first paragraph of a review by R. Penrose of a book by J.C. Polkinghorne “The Quantum World” in The Times Higher Education Supplement, March 23, 1984.

If we were to follow this indication we would be interested in understanding some mathematical properties of semilinear hyperbolic systems with quadratic interaction and the special role played by the four dimensional space-time where we apparently live would probably be linked to Sobolev's embedding theorem, but before playing such a game which accepts too much of the physicists' dogmatical postulates, we should look at a more rational explanation of the mathematical difficulties encountered in the spectroscopic measurements.

Even if we do not know what atoms really are they do appear as tiny obstacles and we have then to face the difficulty of working with at least two scales, a microscopic one and a macroscopic one. In some way the practical goal of quantum physics is to compute corrections in the effective equations satisfied by the macroscopic quantities from a fine and possibly wrong description of what equation the microscopic quantities do satisfy. As mathematicians we should describe a general framework in order to understand more of this question and there are indeed choices to make and difficulties to overcome.

The first obvious choice that we have already made is to work with partial differential equations and not with ordinary differential equations; this can be considered a lesson learned from A. Einstein about the defects of I. Newton's classical approach. There has been much propaganda in recent years for works emphasising finite dimensional effects in partial differential equations and one may indeed be attracted by some of the difficult and interesting mathematical questions which had led H. Poincaré to introduce so many tools and ideas before the development of quantum physics. It was another great mathematician who formalised some of the rules followed by physicists in their quantic games but it is surprising that J. Von Neuman would show that no ordinary differential equation could produce the same results as the rules of quantum mechanics and forget to question the very nature of that set of rules. Certainly if one believed that one should create a game that will generate a sequence of numbers one might be tempted by the mathematical properties of spectra of linear operators. Was then the dogma already well accepted before it became obvious that the spectroscopic experiments were not generating mere lists of numbers? Were mathematicians so impressed by this apparent success of functional analysis? Had there been an intentional effort of propaganda around functional analysis in order to avoid that mathematicians study more relevant partial differential equations of continuum physics in the spirit of what S. Sobolev and J. Leray had already been doing in the 1930s? Certainly, and L. Hörmander [3] is right in pointing at some misconceptions created by L. Schwartz's approach, but some other misconceptions have been propagated by his own approach to partial differential equations. Is there indeed a classical treatise on partial differential equations which does mention these properties of partial differential equations related to the strange effects observed in spectroscopy or more simply which does quote the relation between microscopic and macroscopic levels which is such a crucial question in physics?

The mathematical tool of H-measures which I have introduced [4, 5] is a new step toward a better understanding of these questions and I have chosen the prefix H as a reminder that these objects had arisen naturally in the theory of homogenisation, a term to which I attribute a more general meaning than which is usually given in the rare books related to the subject like those of A. Bensoussan, J.-L. Lions and G. Papanicolaou [6] or of E. Sanchez-Palencia [7] where periodicity assumptions often obscure the methods which I had developed

for more general situations [8], partly in collaboration with F. Murat [9] and as an extension of earlier results of S. Spagnolo [10, 11]. The H-measures added to my previous description of the role of oscillations in partial differential equations [1] that of concentration effects whose importance in continuum physics I had not foreseen before the work of P.-L. Lions, R. DiPerna and A. Majda [12-16]. My initial purpose for introducing H-measures was to derive small amplitude homogenisation theorems [4, 5, 17] in order to explain why some particular formula obtained by physicists [18] was indeed accurate despite the fact that the arguments used in its derivation did not make any sense. I can trace back my intuitive understanding about these objects to some formula for computing an exact quadratic correction term [19] appearing in a model which I had introduced earlier in order to understand some averaging question in hydrodynamics. It was only later that I found a way to use the same H-measures for describing the propagation of oscillations and concentration effects in some partial differential equations [4, 5] obtaining then a quantitative transport property in the form of partial differential equations in  $x$  and  $\xi$  satisfied by the H-measures. I wanted to avoid the standard theory of pseudo-differential operators [3] and construct what I needed for my quadratic microlocal tool of H-measures in order to be able to study partial differential equations of continuum mechanics without making spurious hypotheses of smoothness for the coefficients. However even for those who have devoted a long time reading [3] H-measures may still appear to be natural as they have been introduced independently by P. Gérard [20, 21] although the name of microlocal defect measures which he has chosen for them may reflect a negative attitude inherent in [3].

Of course H-measures are only a step toward the mathematical understanding of these questions of physics which I had sketched at the beginning and there are other pieces of that scientific puzzle which should not be left aside like the apparition of memory effects by homogenisation, which seems the mathematical explanation of what physicists attribute to their strange rules of spontaneous absorption and emission; it must be emphasised that these homogenisation results are obtained without any postulate of a probabilistic nature. There are some more or less classical cases of memory effects induced by homogenisation like viscoelasticity which can be found in Sanchez-Palencia [7] but the effects which I was mentioning are related to hyperbolic situations and have not received much attention apart from my own tentatives [22, 23] and that of Y. Amirat, K. Hamdache and A. Ziani [24, 25] and so a lot remains to be done.

## H-Measures

Contrary to wave front sets which can be attached to general distributions but are mere geometric sets endowed with a negative property of lack of smoothness, H-measures are only defined for sequences of functions converging weakly to zero in  $L^2(\mathbb{R}^N)$  and express in a quantitative way the limit of quadratic quantities, the H-measure being zero in the case of strong convergence in  $L^2(\mathbb{R}^N)$  and, because they are measures on  $\mathbb{R}^N \times S^{N-1}$ , they can see the action of a class of pseudo-differential operators of order zero.

**Definition 1.** An *admissible symbol*  $s$  is a continuous function on  $\mathbb{R}^N \times S^{N-1}$  admitting a decomposition  $s(x, \xi) = \sum_n a_n(\xi) b_n(x)$  with the functions  $a_n$  being

continuous on  $S^{N-1}$ , the functions  $b_n$  being continuous on  $R^N$  and converging to zero at infinity, and such that  $\sum_n \|a_n\| \cdot \|b_n\| < \infty$  where the norms are sup norms. The *standard operator*  $S$  with symbol  $s$  is the continuous operator on  $L^2(R^N)$  defined by  $F(Su)(\xi) = \sum_n a_n(\xi/|\xi|)F(b_n u)(\xi)$  where  $F$  denotes the Fourier transform. A continuous operator  $L$  on  $L^2(R^N)$  is said to have symbol  $s$  if  $L - S$  is a compact operator on  $L^2(R^N)$ .

The only technical point to check is that the commutator  $L_1 L_2 - L_2 L_1$  of two such operators is a compact operator on  $L^2(R^N)$ .

**Proposition 2.** *If  $U^n$  is a sequence converging weakly to zero in  $(L^2(R^N))^p$ , then there is a subsequence and measures  $\mu^{i,j}$  on  $R^N \times S^{N-1}$ ,  $i, j = 1, \dots, p$ , such that for every operators  $L_1, L_2$  with symbols  $s_1, s_2$  the limit of  $L_1(U_i^n)L_2(U_j^n)^*$  is a measure  $v$  on  $R^N$  defined by  $\langle v, \phi \rangle = \langle \mu^{i,j}, \phi s_1 s_2^* \rangle$  for every test function  $\phi$  continuous with compact support in  $R^N$ .*

One immediately finds that  $\mu$  is hermitian nonnegative and has a few other obvious properties, one of them being the following localization principle for H-measures, which is analogous to the information on the wave front sets derived from application of the stationary phase method.

**Proposition 3.** *If a sequence  $U^n$  converges weakly to zero in  $(L^2(R^N))^p$ , corresponds to a H-measure  $\mu$ , and is such that  $\sum_{ij} \partial_i(b_{ij} U_j^n)$  converges strongly to zero in  $H_{loc}^{-1}(R^N)$  where the functions  $b_{ij}$  are continuous, then one has  $\sum_{ij} \xi_i b_{ij} \mu^{jk} = 0$  for  $k = 1, \dots, p$ .*

Before describing the more technical property of propagation let us give a few examples of what was just mentioned.

**Example 4.** Let  $u^n(x) = v(x, x/\varepsilon)$  where  $\varepsilon$  is a sequence converging to zero with  $v$  defined on  $R^N \times R^N$  and  $v(x, y)$  having period 1 in each component  $y_j$ ,  $j = 1, \dots, N$ ; denoting by  $Y$  the unit cube, we assume that  $v$  is continuous in  $x$  with values in  $L^2(Y)$  and decompose  $v$  in Fourier series in  $y$ ,  $v(x, y) = \sum_m v_m(x) e^{2i\pi(m, y)}$ , assuming moreover that  $v_0$  is zero. Under these hypotheses, without extraction of a subsequence, the H-measure  $\mu$  associated to  $u^n$  is defined by  $\langle \mu, \Phi \rangle = \sum_m \int |v_m(x)|^2 \Phi(x, m/|m|) dx$  for every continuous function  $\Phi$  on  $R^N \times S^{N-1}$  with compact support in  $x$ .

**Example 5.** Let  $u^n(x) = \varepsilon^{-N/2} v(x/\varepsilon)$  where  $\varepsilon$  is a sequence converging to zero and  $v$  belongs to  $L^2(R^N)$ . Without extraction of a subsequence the H-measure  $\mu$  associated to  $u^n$  is defined by  $\langle \mu, \Phi \rangle = \int |Fv(\xi)|^2 \Phi(0, \xi/|\xi|) d\xi$  for every continuous function  $\Phi$  on  $R^N \times S^{N-1}$  with compact support in  $x$ .

**Example 6.** Let  $u^n$  be a sequence converging weakly to zero in  $L^2(R^N)$  and corresponding to a H-measure  $\mu$ ; assume moreover that for some continuous

<sup>3</sup> I use L. Schwartz's notations so that  $F(Su)(\xi) = \int s(x, \xi/|\xi|) e^{-2i\pi(x, \xi)} u(x) dx$  for  $u$  smooth with compact support.

<sup>4</sup>  $z^*$  denotes the complex conjugate of  $z$ .

functions  $b_j$ ,  $j = 1, \dots, N$ ,  $\sum_j \partial_j(b_j u^n)$  converges to zero in  $H_{loc}^{-1}(R^N)$  strong. Then  $\mu$  satisfies  $P(x, \xi)\mu = 0$  with  $P$  defined by  $P(x, \xi) = \sum_j b_j(x) \xi_j$ .

**Example 7.** Assume that the sequence  $U^n$  converges weakly to zero in  $(L^2(R^N))^N$  corresponds to a H-measure  $\mu$  and satisfies  $\partial_i U_j^n = \partial_j U_i^n$  for  $i, j = 1, \dots, N$ . Then there exists a scalar nonnegative measure  $v$  on  $R^N \times S^{N-1}$  such that  $\mu^{i,j} = \xi_i \xi_j v$  for  $i, j = 1, \dots, N$ .

**Example 8.** Let  $u^n$  be a sequence converging weakly to zero in  $H^1(R^{N+1})$  and assume that for some continuous functions  $q$  and  $a_{ij}$ ,  $i, j = 1, \dots, N$  independent of  $x_0$ ,  $q\partial_0^2 u^n - \sum_{ij} \partial_i(a_{ij} \partial_j u^n)$  converges to zero in  $H_{loc}^{-1}(R^{N+1})$  strong; assume moreover that  $U^n$  defined by  $U_i^n = \partial_i u^n$  for  $i = 0, \dots, N$ , corresponds to a H-measure  $\mu$ . Then  $\mu^{i,j} = \xi_i \xi_j v$  for  $i, j = 0, \dots, N$  and  $v$  satisfies  $Q(x, \xi)v = 0$  with  $Q$  defined by  $Q(x, \xi) = q(x) \xi_0^2 - \sum_{ij} a_{ij}(x) \xi_i \xi_j$ .

The propagation effects for H-measures are related to the existence of quadratic balance laws and they take the form of partial differential equations in  $(x, \xi)$  satisfied by the H-measures. This is more quantitative than what can be said for wave front sets where it is the complementary property of regularity which is actually propagated. A precise commutation lemma is needed.

**Proposition 9.** Under additional regularity hypotheses, if  $S_1$  and  $S_2$  are the standard operators of symbols  $s_1$  and  $s_2$  then  $\partial_j(S_1 S_2 - S_2 S_1)$  is a continuous operator on  $L^2(R^N)$  with symbol  $\xi_j \{s_1, s_2\}$  where  $\{\cdot, \cdot\}$  denotes the usual Poisson bracket.

In particular if  $s_1 = a(\xi)$  and  $s_2 = b(x)$  then the formula is valid for  $a$  smooth and  $b$  merely of class  $C^1$ , thanks to a result of A. Calderon [26]. In the case of the scalar equation of Example 6 one can obtain then a propagation result under some natural regularity hypotheses.

**Proposition 10.** Let  $u^n$  be a sequence converging weakly to zero in  $L^2(R^N)$  and assume that  $\sum_j b_j \partial_j u^n + c u^n = f^n$  with  $f^n$  converging weakly to zero in  $L^2(R^N)$ , the coefficients  $b_j$  being assumed to be real and of class  $C^1$  while  $c$  is only assumed to be continuous. Assume moreover that  $(u^n, f^n)$  corresponds to a H-measure  $\mu$ . Then  $\mu^{1,1}$  satisfies the following transport equation  $\langle \mu^{1,1}, \{\Phi, P\} - \Phi \operatorname{div} b + 2\Phi \operatorname{Re} c \rangle = \langle 2\operatorname{Re} \mu^{1,2}, \Phi \rangle$  for every function  $\Phi$  of class  $C^1$  in  $(x, \xi)$  with compact support in  $x$ , with  $P$  defined by  $P(x, \xi) = \sum_j b_j(x) \xi_j$ .

In the case of the wave equation of Example 8 one finds a similar result.

**Proposition 11.** Let  $u^n$  be a sequence converging weakly to zero in  $H^1(R^{N+1})$  and assume that  $q\partial_0^2 u^n - \sum_{ij} \partial_i(a_{ij} \partial_j u^n) = f^n$  with  $f^n$  converging weakly to zero in  $L^2(R^N)$ , the coefficient  $q$  being real positive independent of  $x_0$  and of class  $C^1$ , the matrix  $a$  with entries  $a_{ij}$ ,  $i, j = 1, \dots, N$  being hermitian positive independent of  $x_0$  and of class  $C^1$ . Let  $U^n$  be defined by  $U_i^n = \partial_i u^n$  for  $i = 0, \dots, N$ , and assume that  $(U^n, f^n)$  corresponds to a H-measure  $\mu$ . Then  $\mu^{i,j} = \xi_i \xi_j v^{1,1}$  for  $i, j = 0, \dots, N$ , and  $\mu^{i,N+1} = \xi_i v^{1,2}$  for  $i = 0, \dots, N$  and  $v^{1,1}$  satisfies the following transport equation  $\langle v^{1,1}, \{\Phi, Q\} \rangle = \langle 2\operatorname{Re} v^{1,2}, \Phi \rangle$  for every function  $\Phi$  of class  $C^1$  in  $(x, \xi)$  with compact support in  $x$ , with  $Q$  defined by  $Q(x, \xi) = q(x) \xi_0^2 - \sum_{ij} a_{ij}(x) \xi_i \xi_j$ .

One can complete Proposition 10 as in [4] by proving a trace theorem on a noncharacteristic hyperplane and from that deduce a result of change of variables for H-measures under local  $C^1$  diffeomorphism, so that a theory on manifolds could be developed. A more interesting question lies in understanding the effect of semilinearity, that is when  $f^n$  does depend upon  $u^n$  or  $U^n$  in the framework of Proposition 10 or 11. The difficulty lies in the fact that H-measures are intrinsically quadratic objects and do not then predict anything about the limits of trilinear quantities for instance. H-measures do provide an improvement on the method of compensated compactness that I had developed with F. Murat [27, 28] but so far have not provided an alternative approach for quasilinear hyperbolic systems of conservation laws in order to replace the method that I had introduced [28] based on Young measures and compensated compactness, a method which had been successfully applied by R. DiPerna [29, 30, 31, 32]. Ron DiPerna had pointed out many years ago the defects of that old method and the need for a dynamic way of describing oscillations; H-measures is still the best answer to that quest and it is obviously not sufficient. It is important to point out that results like Proposition 11 correspond to the possibility of preparing initial data as a beam concentrated at a point and pointing in some direction and then follow where the energy goes; in the spirit of Example 4 one can also prepare initial data that correspond to a H-measure concentrated at a point in space and charging a countable number of points of the unit sphere and still follow the energy along each of these countably many small beams of light; as wave front sets are closed they cannot even see only a countable dense set of the sphere.

## Other Results

It would be unfair not to point out that P. Gérard has introduced quite interesting variants which I cannot cover here [20, 21]. I cannot either discuss of other applications like the relation with homogenisation [4, 17].

## Conclusion

In conclusion I was quite wrong in that small paragraph spared by the political censors [1] where I associated the study of “propagation of singularities” with classical physics as much better mathematical results connected to the propagation of light are like the above Proposition 11. I was also partly wrong in my previous ideas on quantum physics and oscillations as I had forgotten to include concentration effects in my description. Both these conclusions were learned from the possibilities created by the new mathematical tool of H-measures but a lot remains to be done on the way to a better understanding of physics through increased knowledge of some precise aspects of partial differential equations.

Obviously none of these new results are difficult and they could have been proved a long time ago by any of the best specialists of partial differential equations had they been interested in Science.

*Acknowledgements.* I would not have come here to preach one more time to induce mathematicians to be really interested in Science without the convincing arguments of Michael Crandall; I want to thank him for his continuous support along all these years.

This work is supported by a grant from the National Science Foundation DMS-8803317.

## References

1. Tartar, L.: Etude des oscillations dans les équations aux dérivées partielles non linéaires. In: Trends and Applications of Pure Mathematics to Mechanics. (Lecture Notes in Physics, vol. 195.) Springer, Berlin Heidelberg New York 1984, pp. 384–412
2. Tartar, L.: Moscou sur Yvette, Souvenirs d'un Mathématicien Exclus. 90pp. 1986. Unpublished
3. Hörmander, L.: The analysis of linear partial differential operators I-IV. Springer, Berlin Heidelberg New York 1983–1985
4. Tartar, L.: H-measures, a new approach for studying homogenisation, oscillations and concentration effects in partial differential equations. Proc. Roy. Soc. Edinburgh **115A** (1990) 193–230
5. Tartar, L.: How to describe oscillations of solutions of nonlinear partial differential equations. In: Transactions of the Sixth Army Conference on Applied Mathematics and Computing, ARO Report 89-1, pp. 1133–1141
6. Bensoussan, A., Lions, J.-L., Papanicolaou, G.: Asymptotic analysis for periodic structures. Studies in Mathematics and its Applications, vol. 5. North-Holland, Amsterdam 1978
7. Sanchez-Palencia, E.: Non homogeneous media and vibration theory. (Lecture Notes in Physics, vol. 127.) Springer, Berlin Heidelberg New York 1980
8. Tartar, L.: Cours Peccot. Collège de France, Paris, 1977. Unpublished
9. Murat, F.: H-convergence. Séminaire d'Analyse Fonctionnelle et Numérique, Université d'Alger, 1977–1978
10. Spagnolo, S.: Sul limite delle soluzioni di problemi di Cauchy relativi all' equazione del calore. Ann. Scuola Norm. Sup. Pisa (3) **21** (1967) 657–699
11. Spagnolo, S.: Sulla convergenza di soluzioni di equazioni paraboliche ed ellittiche. Ann. Scuola Norm. Sup. Pisa (3) **22** (1968) 571–597
12. Lions, P.-L.: The concentration compactness principle in the calculus of variations: the locally compact case. Part 1 and 2. Ann. Inst. H. Poincaré, Anal. non Linéaire (1984) 109–145 and no. 4, pp. 223–283
13. Lions, P.-L.: The concentration compactness principle in the calculus of variations: the limit case. Part 1 and 2. Rev. Mat. Iberoamericana **1** (1985) (1) 145–201; (2) 45–121
14. DiPerna, R.J.: Oscillations and concentration in solutions to the equations of mechanics. In: Directions in partial differential equations, M.G. Crandall, P.H. Rabinowitz, R.E.L. Turner (eds.), pp. 43–53. Academic Press, New York 1987
15. DiPerna, R.J., Majda, A.J.: Oscillations and concentration in weak solutions of the incompressible fluid equations. Comm. Math. Phys. **108** (1987) 667–689
16. DiPerna, R.J., Majda, A.J.: Concentration in regularizations for 2-D incompressible flow. Comm. Pure Appl. Math. **40** (1987) 301–345
17. Tartar, L.: H-measures and small amplitude homogenization. In: Random Media and Composites, R.V. Kohn, G. Milton (eds.), pp. 89–99. SIAM, Philadelphia 1989
18. Landau, L.D., Lifschitz, E. M.: Electrodynamics of continuous media. Pergamon Press, 1984
19. Tartar, L.: Remarks on homogenization. In: Homogenization and effective moduli of materials and media, J.L. Ericksen, D. Kinderlehrer, R.V. Kohn, J.-L. Lions (eds.), pp. 228–246. The IMA Volumes in Math. and its Applic., vol. 1. Springer, Berlin Heidelberg New York 1986
20. Gérard, P.: Compacité par compensation et régularité 2-microlocale. Séminaire Equations aux Dérivées Partielles 1988–1989 (Ecole Polytechnique, Palaiseau, exp VI)
21. Gérard, P.: Microlocal defect measures. To appear

22. Tartar, L.: Nonlocal effects induced by homogenization. In: Partial differential equations and the calculus of variations, vol. 2. Essays in Honor of E. DeGiorgi, pp. 925–938. Birkhäuser, 1989
23. Tartar, L.: Memory effects and homogenization. To appear in Arch. Rational Mech. Anal.
24. Amirat, Y., Hamdache, K., Ziani A.: Homogénéisation d'équations hyperboliques du premier ordre et application aux écoulements miscibles en milieu poreux. Ann. Inst. H. Poincaré, Anal. non Linéaire (5) **6** (1989) 397–417
25. Amirat, Y., Hamdache, K., Ziani A.: Etude d'une équation de transport à mémoire. C. R. Acad. Sci. Paris, Sér. I Math. **311** (1990) 685–688
26. Calderon, A.P.: Commutators of singular integral operators. Proc. Nat. Acad. Sci. **53** (1978) 1092–1099
27. Murat, F.: Compacité par compensation. Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4) **5** (1978) 489–507
28. Tartar, L.: Compensated compactness and applications to partial differential equations. In: Nonlinear analysis and mechanics. Heriot-Watt Symposium IV, R.J. Knops (ed.). Research Notes in Mathematics, vol. 39, pp. 136–212. Pitman, London 1979
29. DiPerna, R.J.: Convergence of approximate solutions to conservation laws. Arch. Rat. Mech. Anal. **82** (1983) 27–70
30. DiPerna, R.J.: Convergence of the viscosity method for isentropic gas dynamics. Comm. Math. Phys. **91** (1983) 1–30
31. Tartar, L.: The compensated compactness method applied to systems of conservation laws. In: Systems of Nonlinear Partial Differential Equations, J.M. Ball (ed.). NATO ASI Series C 111, pp. 263–285. Reidel, New York 1983
32. Tartar, L.: Discontinuities and oscillations. In: Directions in partial differential equations, M.G. Crandall, P.H. Rabinowitz, R.E.L. Turner (eds.), pp. 211–233. Academic Press 1987



# Microlocal Analysis in Spectral and Scattering Theory and Index Theory

Michael E. Taylor

Department of Mathematics, University of North Carolina  
Chapel Hill, NC 27599-3250, USA

## 1. Introduction

Microlocal analysis has the dual foundation of the theory of singular integral operators, which arose originally to treat elliptic PDE, and geometrical optics, which arose to describe solutions to wave equations. It was forged into a powerful general tool for linear PDE a little over 20 years ago, through the efforts of at least five groups, including Calderon and Zygmund and their students; Kohn and Nirenberg; Sato and his collaborators; Maslov and Egorov and their collaborators; and Hörmander.

The basic object of microlocal analysis is the Fourier integral operator, which can be written as an oscillatory integral

$$Au(x) = \int a(x, y, \theta) e^{i\psi(x, y, \theta)} u(y) dy d\theta, \quad (1.1)$$

where  $\psi$  is a real valued phase function, typically homogeneous of degree 1 in  $\theta$ , and the amplitude  $a(x, y, \theta)$  belongs to a symbol class, typically characterized by estimates on its derivatives, e.g.,

$$|D_y^\alpha D_x^\beta D_\theta^\gamma a(x, y, \theta)| \leq C_{\alpha\beta\gamma} \langle \theta \rangle^{m-\varrho|\alpha|+\delta|\beta|+\delta|\gamma|} \quad (1.2)$$

to define  $S_{\varrho,\delta}^m$ . If  $a$  is asymptotic to a sum of terms homogeneous of degree  $m-j$  in  $\theta$ ,  $j = 0, 1, 2, \dots$ , we say  $a \in S^m$ . An important special case of (1.1) is

$$Au(x) = \int a(x, \xi) e^{i\varphi(x, \xi)} \hat{u}(\xi) d\xi, \quad (1.3)$$

where  $\hat{u}$  is the Fourier transform of  $u$ ; the phase function is  $\psi = \varphi(x, \xi) - y \cdot \xi$ . One imposes a nondegeneracy condition implying that  $(\nabla_\xi \varphi, \xi) \mapsto (x, \nabla_x \varphi)$  is locally well defined; this is the canonical transformation associated to (1.3). Under appropriate conditions on  $\varphi$  in (1.1), one obtains a Fourier integral operator associated to a more general sort of canonical relation. The identity canonical transformation arises from the phase function  $x \cdot \xi - y \cdot \xi$ , for which (1.3) specializes to the formula

$$Au(x) = \int a(x, \xi) e^{ix \cdot \xi} \hat{u}(\xi) d\xi \quad (1.4)$$

for a pseudodifferential operator. If  $a \in S_{\rho,\delta}^m$ , we say the operator  $A$  in (1.4) belongs to  $OPS_{\rho,\delta}^m$ .

Other sorts of pseudodifferential operators and Fourier integral operators have also arisen over the past 20 years, as tools in the study of various problems in PDE. Some of them, useful for certain problems in spectral theory and in scattering theory, will be described below.

## 2. Diffraction Effects

Operators of the form (1.3) describe propagation of waves away from a boundary, and also transversal reflection of waves off a boundary. Following earlier formal developments, Melrose [Me1] and Taylor [T1] made rigorous analyses of propagation of singularities along rays hitting a boundary at grazing incidence, producing Fourier-Airy operators of the form

$$L(F) = \int [gA(\zeta) + ihA'(\zeta)] A(\zeta_0)^{-1} e^{i\theta} \hat{F}(\zeta) d\zeta, \quad (2.1)$$

in which  $A = A_{\pm}$  is an Airy function,  $A_{\pm}(z) = Ai(e^{\mp 2\pi i/3} z)$ ,  $\theta(x, \xi), \zeta(x, \xi)$  is a pair of phase functions, and  $g(x, \xi), h(x, \xi)$  is a pair of amplitudes. The phase functions satisfy a system of eikonal equations and the amplitudes satisfy a system of transport equations. The parametrix (2.1) applies for example to solutions to the wave equation  $u_{tt} - \Delta u = 0$  on  $\Omega = \mathbb{R} \times \mathcal{O}$  where  $\mathcal{O} = \mathbb{R}^n \setminus K$ ,  $K$  smooth and strictly convex, with lateral boundary  $\partial\Omega = \mathbb{R} \times \partial K$ . The Dirichlet condition can be formulated as  $u = f$  on  $\partial\Omega$  with  $f \in \mathscr{C}'(\partial\Omega)$ , and we require  $u = 0$  for  $t < 0$ . Then  $L(F)|_{\partial\Omega}$  is given by an elliptic Fourier integral operator  $J$ , acting on  $F$ , and the parametrix is given by (2.1) with  $F = J^{-1}(f)$ .

An important associated operator is the Neumann operator  $N$ , defined by  $Nf = \partial u / \partial v|_{\partial\Omega}$  where  $u$  solves the Dirichlet problem with data  $f$ . From (2.1) one obtains the fundamental formula

$$N = J(A_1 \Phi + B_1) J^{-1} \quad (2.2)$$

for a certain elliptic  $A_1 \in OPS^{\frac{1}{2}}$ ,  $B_1 \in OPS^0$ . Here  $\Phi$  is defined by  $(\Phi f)^{\wedge}(\xi) = \Phi(\zeta_0)\hat{f}(\xi)$ , where  $\zeta_0 = |\xi|^{-\frac{1}{2}}\xi_n$  appears in (2.1), and  $\Phi(z) = A'(z)/A(z)$  is an Airy quotient. The operator  $A_1 \Phi + B_1$  belongs to  $OPS_{1/3,0}^1$ ; some of its operator properties are discussed in §3 below.

In [MeT1]–[MeT2] this parametrix construction was applied to produce some detailed results on scattered high frequency waves, relating to the behavior of the outgoing solution  $v(x, \lambda)$  to

$$(\Delta + \lambda^2)v = 0 \text{ on } \mathbb{R}^n \setminus K, \quad v = e^{-i\lambda x \cdot \omega} \text{ on } \partial K. \quad (2.3)$$

This included a corrected Kirchhoff approximation, describing the behavior of  $\partial_v v(x, \lambda)$  for  $x \in \partial K$ . We obtained an expansion of the form  $K(\omega, x, \lambda)e^{-i\lambda x \cdot \omega}$  with

$$K(\omega, x, \lambda) \sim \sum_{j,k \geq 0} \lambda^{\frac{2}{3}-k-\frac{2}{3}j} b_{jk}(\omega, x) \Psi^{(j)}(\lambda^{\frac{1}{3}} Z), \quad (2.4)$$

$Z$  being a smooth function vanishing at the shadow boundary, determined in a subtle fashion by the symplectic geometry (in particular, *not* generally equal to  $v \cdot \omega$ ). The function  $\Psi$  is given by

$$\Psi(\tau) = e^{-i\tau^3/3} \int A(s)^{-1} e^{-is\tau} ds; \quad (2.5)$$

it is rapidly decreasing as  $\tau \rightarrow -\infty$  and satisfies

$$\Psi(\tau) \sim \sum_{j \geq 0} \alpha_j \tau^{1-3j}, \quad \tau \rightarrow +\infty. \quad (2.6)$$

This appears to sharpen and extend some formulas of Fok. It demonstrates rigorously the often divined transitional layer of width  $\sim \lambda^{-\frac{1}{3}}$  about the shadow boundary. Work on the corrected Kirchhoff approximation for solutions to Maxwell's equations was done in Yingst [Yi].

The asymptotic behavior of  $v(x, \lambda)$  away from  $\partial K$  exhibits *two* transitional regions near the shadow boundary, one of width  $\sim \lambda^{-\frac{1}{3}}$  and another of width  $\sim \lambda^{-\frac{1}{2}}$ . Near the shadow boundary, it was shown in [MeT2] that

$$v(x, \lambda) \sim \int b(x, \zeta, \lambda) \left( \frac{A_+}{A_-} \right) (\lambda^{\frac{2}{3}} \zeta) e^{i\lambda\varphi(x, \zeta)} d\zeta \quad (2.7)$$

where  $b$  has compact support in  $\zeta$  and an asymptotic expansion in powers of  $\lambda$ . This has a further expansion of the form

$$\lambda^\alpha [e^{i\lambda p_0} p_0(x, \lambda) + e^{i\lambda p_1} p_1(x, \lambda) + e^{i\lambda p_2} p_2(x, \lambda)]. \quad (2.8)$$

The first term captures the direct wave, the third term the reflected wave, and the middle term represents a diffraction effect confined very close to the shadow boundary. The amplitudes  $p_j(x, \lambda)$  have the following form:

$$p_0(x, \lambda) = a_0(x, \lambda^{\frac{1}{2}} s_0(x), \lambda^{\frac{1}{2}}) \quad (2.9)$$

where here and below  $s_j(x)$  are smooth functions vanishing simply on the shadow boundary.  $a_0(x, \tau, \mu)$  is a symbol of product type in  $\tau, \mu$ .

$$p_2(x, \lambda) = a_2(x, \lambda^{\frac{1}{3}} s_2(x), \lambda^{\frac{1}{3}}), \quad (2.10)$$

and  $a_2(x, \tau, \mu)$  is also a symbol of product type. The most subtle term is

$$p_1(x, \lambda) = a_1(x, \lambda^{\frac{1}{3}} s_1(x), \lambda^{-\frac{1}{6}}) \quad (2.11)$$

where  $a_1(x, \tau, \mu)$  is  $C^\infty$  away from  $(\tau, \mu) = (0, 0)$ , rapidly decreasing as  $\tau \rightarrow \infty$ , and has a conormal singularity at  $(0, 0)$ , which can be approached as  $\lambda \rightarrow \infty$  due to the negative exponent on  $\lambda$  in the last argument of  $a_1$ .

Zworski [Zw1] has extended the analysis of the asymptotic behavior of  $v(x, \lambda)$  to be uniformly valid as  $x$  approaches  $\partial K$ . In [Zw2] a rigorous study of the shift in the shadow boundary predicted by Keller and Rubinow is made; the shift is asymptotic to  $C_0 \lambda^{-\frac{2}{3}}$ , well inside the smaller transitional region, hence in a region where the asymptotic expansion (2.8) simplifies.

In [MeT1] there is also a uniform analysis of the near peak scattering amplitude, extending earlier work of [MjT] and [Me3] on the on peak behavior. The expansion is somewhat like (2.8). Again there are two transitional regions and a subtle term connecting them. This time the widths of the transition regions are  $\sim \lambda^{-\frac{1}{2}}$  and  $\sim \lambda^{-1}$ .

A key ingredient in these analyses was a normal form for Fourier integral operators with folding canonical relations, namely

$$P_1 \mathcal{A}i + P_2 \mathcal{A}'i, \quad P_1 \in OPS^m, \quad P_2 \in OPS^{m-\frac{1}{2}}, \quad (2.12)$$

where  $\mathcal{A}i$  is Fourier multiplication by  $Ai(\zeta_0)$  and  $\mathcal{A}'i$  is similarly defined. Using this, Farris [Fa] found a microlocal model for solutions to boundary problems with grazing rays, namely

$$K_1(\mathcal{A}_+/\mathcal{A}_-)K_2 \quad (2.13)$$

where  $K_j$  are elliptic Fourier integral operators and  $(\mathcal{A}_+/\mathcal{A}_-)$  is Fourier multiplication by  $A_+(\zeta_0)/A_-(\zeta_0)$ .

Gliding ray problems arise for example if one considers the wave equation  $u_{tt} - \Delta u = 0$  on  $\mathbb{R} \times K$  rather than  $\mathbb{R} \times \mathcal{O}$ , where as above  $K \subset \mathbb{R}^n$  has smooth strictly convex boundary. In such a case, one has a parametrix like (2.1) with  $A_\pm(\zeta)$  replaced by  $Ai(\zeta)$ . Since  $Ai(z)$  has real zeros, it is convenient to make an almost analytic continuation of  $\zeta(x, \xi)$  and evaluate it at  $\xi_n + iT$  for some  $T \neq 0$  rather than at  $\xi_n$ . Then one has instead of (2.2) a formula for the Neumann operator involving  $\Phi i$ , Fourier multiplication by  $Ai'(\zeta_0)/Ai(\zeta_0)$ . It is of great help that, thanks to Melrose's work on equivalence of glancing hypersurfaces, one can arrange that  $\zeta = \zeta_0$  on the boundary,  $\zeta_0(\xi) = |\xi|^{-\frac{1}{2}}\xi_n$ . Not having this introduces complications, which however were tackled in [Es1]. One needs to resort to energy estimates rather than use symbol calculus, and hence produces a less explicit sort of parametrix.

Having the replacement for (2.1)–(2.2) described above for gliding ray problems does not end one's job here. The operator  $\Phi i$  is a rather complicated operator, a singular sort of Fourier integral operator with an infinite number of canonical transformations accumulating along the leaf relation of the characteristic variety for the boundary. One key method for taming  $\Phi i$  is described in the next section.

### 3. Airy Operator Calculus

Solving boundary problems other than the Dirichlet problem for wave equations with grazing rays involves, in addition to (2.1), the solution to an equation of the form

$$(A\Phi + B)f = g, \quad (3.1)$$

with  $A \in OPS^{\frac{2}{3}}$ ,  $B \in OPS^1$ . If  $B$  is elliptic, one has an elliptic operator in  $OPS_{\frac{1}{3}, 0}^1$ , while if  $A$  is elliptic and the symbol of  $B$  vanishes on  $\{\zeta_0 = 0\}$ , one has a hypoelliptic operator. There is a special calculus of Airy operators, extending that of [Me2], sketched in [T5] and treated in detail in [MeT3], which provides more analytical detail on compositions and parametrices of such operators than

the  $S_{1/3,0}^m$  calculus, based on some remarkable identities which follow from constructing the Neumann operator (2.2) using different choices of solutions  $g, h$  to the transport equations for the amplitudes in (2.1).

The full advantages of this approach are apparent when one treats the analogues that arise in problems with gliding rays, i.e., composition and parametrices for operators like

$$A\Phi i + B, \quad (3.2)$$

with  $A$  and  $B$  as in (3.1). These are highly singular variants of Fourier integral operators, with a rather subtle operator calculus. Nevertheless, if  $B$  is elliptic, the operator (3.2) is shown to have a parametrix in the class  $\mathcal{A}_\sigma^{i^{-1}, \pm}$ , where by definition  $\mathcal{A}_\sigma^{i^m, \pm}$  consists of operators with asymptotic expansions of the form

$$T \sim B + \sum_{j \geq 0} A_j \Phi i C_j, \quad B \in OPS^m, \quad A_j \in OPS^{m_j}, \quad C_j \in OPS^0, \quad (3.3)$$

where  $m_j + \frac{1}{3} = m - \ell_j, \ell_j \geq 0$  is an integer,  $\ell_j \rightarrow \infty$  as  $j \rightarrow \infty$ . The sign  $+$  or  $-$  reflects the choice of sign of  $T$ , with  $\zeta_0$  evaluated at  $\xi_n + iT$ . If  $A$  is elliptic and  $B$  has vanishing symbol on  $\{\zeta_0 = 0\}$ , the operator (3.2) has a parametrix of the form

$$(C\Phi i^{-1} + D)R(1 + E\Phi i^{-1}R)^{-1} \quad (3.4)$$

with  $C \in OPS^{-1}, D \in OPS^{-\frac{4}{3}}, R \in OPS^{\frac{1}{3}}$  elliptic,  $E \in OPS^{-1}$ . Then

$$E\Phi i^{-1}R : H^s \longrightarrow H^{s+\frac{1}{3}}, \quad (3.5)$$

so the Neumann series for the last factor is asymptotic. This affords a symbolic construction of parametrices for numerous boundary problems involving gliding rays.

There are alternative approaches to many boundary problems, both in grazing and gliding cases, that involve applying energy estimates to such equations as (3.1) or its analogue using the operator (3.2). Information so obtained is less explicit than by a symbolic parametrix construction, but can be useful all the same. Energy estimate approaches also have a flexibility to apply to cases not amenable to explicit constructions. Eskin [Es2] has proposed a number of techniques along those lines. In [T6] the Fefferman-Phong inequality was brought to bear on a number of equations of the form (3.1).

## 4. Functional Calculus

Let  $A$  be a self adjoint elliptic operator in  $OPS^1$ . It is convenient to analyze many functions of  $A$  via

$$f(A) = \int \hat{f}(t) e^{itA} dt, \quad (4.1)$$

since  $e^{itA}$  is a group of Fourier integral operators, of the form (1.3) for  $|t|$  small. If  $f$  belongs to a symbol space,  $\hat{f}(t)$  is singular only at the origin, and is rapidly decreasing as  $|t| \rightarrow \infty$ . One can write  $\hat{f}(t) = \hat{f}_1(t) + \hat{f}_2(t)$  where  $\hat{f}_1$  has support near

$t = 0$  and  $\hat{f}_2 \in \mathcal{S}(\mathbb{R})$ . We can use (1.3) to analyze  $f_1(A)$  as a pseudodifferential operator. If  $A$  acts on functions on a compact manifold  $M$ , then  $f_2(A) \in OPS^{-\infty}$  is for many purposes negligible. If  $M$  is not compact,  $A = (c - \Delta)^{\frac{1}{2}}$ , (with  $c - \Delta \geq 0$ ), and  $f$  is even, one can replace (4.1) by

$$f(A) = \int \hat{f}(t) \cos tA \, dt, \quad (4.2)$$

and exploit finite propagation speed to get useful information on  $f_2(A)$ .

Using (4.1) to evaluate the trace of  $f(A)$  and hence give sharp results on the spectral asymptotics of  $A$  was one of the early spectacular applications of Fourier integral operators. In Hörmander's paper [Ho1], a key point is to analyze the trace of (4.1) for  $f(\lambda) = f_\tau(\lambda) = f(\lambda - \tau)$  as  $\tau \rightarrow \infty$ , given  $\hat{f} \in C_0^\infty(\mathbb{R})$  having small support. [DG] extended this to any  $\hat{f} \in C_0^\infty(\mathbb{R})$ . In [T2] and Chapter 12 of [T3] a number of applications of the fact that  $f(A) \in OPS_{1,0}^n$  when  $f \in S_1^n(\mathbb{R})$  were discussed. It was also shown how functional calculus applied to the Laplace operator on  $S^2$  led to a clean treatment of the problem of scattering of waves in  $\mathbb{R}^3$  by a sphere, a special case of the general problem discussed in §1 which has been treated by many authors, generally with a heavier dependence on special function theory. About the same time, Colin de Verdiere [CV] discussed a similar theory of functions of pseudodifferential operators, emphasizing applications to semiclassical asymptotics for spectra of Schrödinger operators.

The method of separation of variables, combined with a harmonic analysis on the base using the techniques described above, were used by [CT] to give a detailed analysis of diffraction of waves by a cone.

The formula (4.2) was exploited in [CGT] to produce fine estimates on the heat kernel on a variety of complete Riemannian manifolds. Of particular interest were those for which the Ricci tensor was bounded from below; also stronger hypotheses, such as  $C^\infty$ -bounded geometry, yielded stronger results, such as  $L^p$ -boundedness of (4.2) for  $f \in S_1^0(\mathbb{R})$  holomorphic in a strip of width related to the volume growth of  $M$ . In [CGT] we used this to recover known results on  $L^p$ -boundedness when  $M$  is a symmetric space of rank 1. In [DST] a study was made of the  $L^p$ -spectrum of the Laplace operator on geometrically finite quotients of hyperbolic space. In [T8] the argument of [CGT] was honed and yielded results on general symmetric spaces of noncompact type, sharp enough to establish conjectures on what is precisely the  $L^p$ -spectrum of the Laplace operator. In such a case, as is well known, the  $L^2$ -spectrum of  $-\Delta$  is  $[\|\varrho\|^2, \infty)$ . Then the  $L^p$ -spectrum is shown to be precisely  $\{|\varrho|^2 + z^2 : |\text{Im}z| \leq |\frac{2}{p} - 1| \cdot |\varrho|\}$ . Furthermore, with  $H = -\Delta - |\varrho|^2 \geq 0$ ,  $L = H^{\frac{1}{2}}$ ,  $f(L)$  is bounded on  $L^p$  provided  $1 < p < \infty$  and  $f$  is holomorphic on a strip  $|\text{Im}z| \leq |\frac{2}{p} - 1| \cdot |\varrho|$  and is a symbol of order 0 there. Recently, J. Anker [An], [An2] has further extended these arguments.

## 5. Semiclassical Asymptotics and Gauge Fields

Let  $M$  be a Riemannian manifold,  $G$  a compact Lie group,  $P \rightarrow M$  a principal  $G$ -bundle. A gauge field is defined by a connection on  $P$ . Corresponding to each irreducible unitary representation  $\lambda$  of  $G$  is a Hermitian vector bundle  $E_\lambda \rightarrow M$  with connection,  $\nabla_\lambda$ , and a Hamiltonian operator  $H_\lambda^0 = \nabla_\lambda^* \nabla_\lambda$ . If a scalar potential  $V$  is also given, a semiclassical analysis of the resulting quantum system involves study of the spectrum of

$$H_\lambda = \hbar^2 H_\lambda^0 + V, \quad \hbar = |\lambda + \delta|^{-1}, \quad (5.1)$$

as  $\hbar \rightarrow 0$ . Thus  $\lambda \rightarrow \infty$  in a Weyl chamber. This problem was studied in [ST1]–[ST2]. In particular, given  $f \in \mathcal{S}(\mathbb{R})$ , an asymptotic expansion was made of  $\text{Tr } f(H_\lambda)$ ; it was shown that

$$\text{Tr } f(H_\lambda) = \langle \kappa, \chi_\lambda \rangle = d_\lambda \beta(\lambda + \delta) \quad (5.2)$$

for a Weyl group invariant symbol  $\beta$ , defined by the spectrum of a certain central distribution  $\kappa$  on  $G$  arising from the operator  $f(-A^{-1}L) \in OPS^0(P)$ , where

$$L = \Delta + (V - 1)\Delta_G^P - |\delta|^2 V, \quad A = -\Delta_G^P + |\delta|^2, \quad (5.3)$$

$\Delta$  being the Laplace operator on  $P$  and  $\Delta_G^P$  the vertical Laplacian. One arranges that  $V > 1$ . Methods of §4 were used to analyze  $f(-A^{-1}L)$ . We have  $\kappa = \text{Tr}_G f(-A^{-1}L)$ , where the ‘ $G$ -trace’ of an operator on  $C^\infty(P)$  with Schwartz kernel  $K(p, q)$  is defined to be

$$\kappa(g) = \int_P K(p \cdot g, p) dV(p). \quad (5.4)$$

In [GU], using different techniques, a different sort of asymptotic study of the spectrum of (5.1) was made, as  $\lambda \rightarrow \infty$  along a ray. In [TU] a synthesis of these results has been achieved. We study the asymptotic behavior of

$$\text{Tr } f(\hbar^{-1} H_\lambda^{-\frac{1}{2}} (H_\lambda - c)), \quad \hbar = |\lambda + \delta|^{-1} \quad (5.5)$$

as  $\hbar \rightarrow 0$ , for a given  $c \in \mathbb{R}$ , a regular value of  $V$ , given  $\hat{f} \in C_0^\infty(\mathbb{R})$ . This is analyzed in terms of the  $G$ -trace of  $f(Q)$ , where

$$Q = (-L)^{-\frac{1}{2}} (-L - cA) \in OPS^1(P). \quad (5.6)$$

Now  $Q$  is typically not elliptic, but it is an operator of principal type, with real principal symbol. A modification of the method of §4 shows that, for  $\hat{f} \in C_0^\infty(\mathbb{R})$ ,  $f(Q)$  is a Fourier integral operator on  $C^\infty(P)$  with a canonical relation given by the leaf relation on the characteristic variety of  $Q$ . Then  $\text{Tr}_G f(Q)$  is under reasonable geometrical conditions a central Lagrangian distribution on  $G$ , whose nature reflects the classical dynamics of the gauge and scalar field on the Wong bundle. Harmonic analysis of  $\text{Tr}_G f(Q)$  involves restriction to a maximal torus  $T$ , an operation for which clean intersection hypotheses often fail, if  $G$  has rank  $\geq 2$ . In the simplest cases there can arise distributions on  $T$  associated to transversally intersecting Lagrangians, thus giving rise to nonclassical asymptotics for (5.5).

## 6. Operator $K$ -Theory

Atiyah proposed  $K$ -homology as an abstract setting for index theory, and this was developed by Brown, Douglas, and Fillmore and by Kasparov and others; see [BII]. In [BDT] an investigation was made of cycles defined by elliptic differential operators. It was shown that every closed extension of a first order elliptic differential operator  $D$ , acting on sections of a vector bundle over a smooth Riemannian manifold  $M$ , not necessarily complete, defines a class  $[D]$  in  $KK(C_0(M), \mathbb{C})$ , independent of the choice of closed extension. The proof involves a use of (4.1) and finite propagation speed. In case  $M$  is contained in  $\bar{M}$ , compact with boundary, consideration of the boundary map  $\delta : KK(C_0(M), \mathbb{C}) \rightarrow K_1(\partial M)$  applied to cycles coming from different closed extensions leads to identities in  $K$ -homology.

For example, if  $M = \Omega$  is a pseudoconvex manifold with the property that the  $\bar{\partial}$ -Neumann problem has compact resolvent on  $(0, p)$ -forms, for  $p \neq 0$ , then taking two natural closed extensions of  $D = \bar{\partial} + \bar{\partial}^*$  and applying the boundary map produces the identity

$$[D_{\partial\Omega}] = [\tau_\Omega] \text{ in } K_1(\partial\Omega) \quad (6.1)$$

where  $D_{\partial\Omega}$  is the Dirac operator on  $\partial\Omega$ , with its natural  $spin^c$ -structure, and  $[\tau_\Omega]$  is the Toeplitz extension. This identity refines and generalizes Boutet de Monvel's index theorem for elliptic Toeplitz operators. It also leads to a number of other identities in  $K$ -homology.

Identities in  $K_1(M)$ , including (6.1), can be used together with the Bott map to produce identities in  $K_0(M)$ .

In [T9] an examination was made of ways in which differential operators and pseudodifferential operators define elements of  $K^j(\Psi^0(M))$ , where  $\Psi^0(M)$  is the  $C^*$ -algebra which is the  $L^2$ -operator norm closure of  $OPS^0(M)$ , in case  $M$  is compact. It was shown that in some cases (6.1) has a further refinement as an identity in  $K^1(\Psi^0(\partial\Omega))$ . I suspect that in general the two elements differ by a quantity which can be described in terms of

$$\tau : K^0(M) \rightarrow K^1(\Psi^0(M)), \quad (6.2)$$

which arose in [T9] to describe an obstruction for an elliptic pseudodifferential operator acting on sections of a vector bundle to define an element of  $K^0(\Psi^0(M))$ .

The  $K$ -theory of  $\Psi^0(M)$  can be thought of as a microlocal version of the  $K$ -theory of  $M$ . Extra structure arises, partly due to the richer structure of ideals; in particular, if  $A \subset T^*M \setminus 0$  is a closed conic set, one can form  $\Psi_A^0(M)$ , the closure of the set of elements of  $OPS^0(M)$  whose principal symbols vanish on  $A$ . For example, if  $M$  has a contact structure, one can let  $A$  be the contact line bundle  $(\setminus 0)$ ; then  $\Psi^0/\Psi_A^0 \approx C(M) \oplus C(M)$ . To give one example of a natural cycle arising in this context, if  $M = \partial\Omega$  is the boundary of a strongly pseudoconvex domain in  $\mathbb{C}^2$ ,  $\bar{\partial}_b$  defines an element of  $KK(\Psi_A^0(M), \mathbb{C})$ . It is of interest to compute the (co)boundary map  $\delta : KK(\Psi_A^0, \mathbb{C}) \rightarrow K^1(\Psi^0/\Psi_A^0)$ , isomorphic to  $K_1(M) \oplus K_1(M)$  in this last case. It is noted in [T9] that, with  $D_M$  as in (6.1),

$$\delta[\bar{\partial}_b] = ([D_M], -[D_M]). \quad (6.3)$$

## References

- [An] Anker, J.:  $L^p$  Fourier multipliers on Riemannian symmetric spaces of the non-compact type. Preprint
- [An2] Anker, J.: Handling the inverse spherical Fourier transform. Preprint
- [BDT] Baum, P., Douglas, R., Taylor, M.: Cycles and relative cycles in analytic K-homology. *J. Diff. Geom.* **30** (1989) 761–804
- [BlJ] Blackadar, B.: K-theory for operator algebras. Springer, Berlin Heidelberg New York 1986
- [BdM] Boutet de Monvel, L.: On the index of Toeplitz operators of several complex variables. *Invent. math.* **50** (1979) 249–272
- [CGT] Cheeger, J., Gromov, M., Taylor, M.: Finite propagation speed, kernel estimates for functions of the Laplace operator, and the geometry of complete Riemannian manifolds. *J. Diff. Geom.* **17** (1982) 15–53
- [ChT] Cheeger, J., Taylor, M.: Diffraction of waves by conical singularities. *Comm. Pure Appl. Math.* **35** (1982) 275–331, 487–529
- [CV] Colin de Verdiere, Y.: Spectre conjoint d'opérateurs pseudodifferentiels qui commutent, I: Le cas non intégrable. *Duke Math. J.* **46** (1979) 169–182. II: Le cas intégrable. *Math. Z.* **171** (1980) 51–73
- [Cor] Cordes, H.O.: Elliptic pseudo-differential operators – An abstract theory. (Lecture Notes in Mathematics, vol. 756.) Springer, Berlin Heidelberg New York 1979
- [DST] Davies, E.B., Simon, B., Taylor, M.:  $L^p$  spectral theory of Kleinian groups. *J. Funct. Anal.* **78** (1988) 116–136
- [DG] Duistermaat, J.J., Guillemin, V.: The spectrum of positive elliptic operators and periodic bicharacteristics. *Invent. math.* **29** (1975) 39–79
- [Es1] Eskin, G.: Parametrix and propagation of singularities for the interior mixed hyperbolic problem. *J. Anal. Math.* **32** (1977) 17–62
- [Es2] Eskin, G.: General initial-boundary problems for second order hyperbolic equations. In: *Singularities in boundary value problems*. D. Reidel, Boston 1981, pp. 19–54
- [Fa] Farris, M.: Egorov's theorem on a manifold with diffractive boundary, *Comm. PDE* **6** (1981) 651–688
- [GU] Guillemin, V., Uribe, A.: Reduction, the trace formula, and semiclassical asymptotics. *Proc. Nat. Acad. Sci. USA* **84** (1987) 7799–7801
- [HR] Hellfer, B., Robert, D.: Comportement semi-classique du spectre des Hamiltoniens quantiques périodiques. *Ann. Inst. Fourier* **31** (1981) 169–223
- [Ho1] Hörmander, L.: The spectral function of an elliptic operator. *Acta Math.* **121** (1968) 193–218
- [Ho2] Hörmander, L.: The analysis of linear partial differential operators, vols. 3 and 4. Springer, Berlin Heidelberg New York 1985
- [Ki] Kirchhoff, G.: Vorlesungen über Math. Physik (Optik). Leipzig, 1891.
- [MjT] Majda, A., Taylor, M.: The asymptotic behavior of the diffraction peak in classical scattering. *Comm. Pure Appl. Math.* **30** (1977) 639–669
- [Me1] Melrose, R.: Microlocal parametrices for diffractive boundary value problems. *Duke. Math. J.* **42** (1975) 605–635
- [Me2] Melrose, R.: Airy operators. *Comm. PDE* **3** (1978) 1–76
- [Me3] Melrose, R.: Forward scattering by a convex obstacle. *Comm. Pure Appl. Math.* **33** (1980) 461–499
- [MeS] Melrose, R., Sjöstrand, J.: Singularities of boundary value problems. I.: *Comm. Pure Appl. Math.* **31** (1978) 593–617, II: *Comm. Pure Appl. Math.* **35** (1982) 129–168

- [MeT1] Melrose, R., Taylor, M.: Near peak scattering and the corrected Kirchhoff approximation for convex bodies. *Adv. Math.* **55** (1985) 242–315
- [MeT2] Melrose, R., Taylor, M.: The radiation pattern of a diffracted wave near the shadow boundary. *Comm. PDE* **11** (1986) 599–672
- [MeT3] Melrose, R., Taylor, M.: Boundary problems for wave equations with grazing and gliding rays. To appear
- [ST1] Schrader, R., Taylor, M.: Small  $\hbar$  asymptotics for quantum partition functions associated to particles in external Yang-Mills potentials. *Commun. Math. Phys.* **92** (1984) 555–594
- [ST2] Schrader, R., Taylor, M.: Semiclassical asymptotics, gauge fields, and quantum chaos. *J. Funct. Anal.* **83** (1989) 258–316
- [Str] Strichartz, R.: A functional calculus for elliptic pseudodifferential operators. *Amer. J. Math.* **94** (1972) 711–722
- [T1] Taylor, M.: Grazing rays and reflection of singularities of solutions to wave equations. *Comm. Pure Appl. Math.* **29** (1976) 1–38, 463–481
- [T2] Taylor, M.: Fourier integral operators and harmonic analysis on compact manifolds. *AMS Proc. Symp. Pure Math.* **35** (2) (1979) 115–136
- [T3] Taylor, M.: Pseudodifferential operators. Princeton Univ. Press, 1981
- [T4] Taylor, M.: Diffraction effects in the scattering of waves. In: *Singularities in boundary value problems*, D. Reidel (ed.). Boston 1981, pp. 271–316
- [T5] Taylor, M.: Airy operator calculus. *Contemp. Math. AMS* (1984) 169–192
- [T6] Taylor, M.: Fefferman-Phong inequalities in diffraction theory. *AMS Proc. Symp. Pure Math.* **43** (1984) 261–300
- [T7] Taylor, M.: Noncommutative microlocal analysis. *Memoirs AMS* **313** (1984)
- [T8] Taylor, M.:  $L^p$ -estimates on functions of the Laplace operator. *Duke Math. J.* **58** (1989) 773–793
- [T9] Taylor, M.: Pseudodifferential operators and K-homology. I: *AMS Proc. Symp. Pure Math.* **51** (1) 561–583; II: *Contemp. Math.* **106** (1990) 245–269
- [T10] Taylor, M.: Pseudodifferential operators and nonlinear PDE. *Progress in Math.* Birkhäuser, 1991 (to appear)
- [TU] Taylor, M., Uribe, A.: Semiclassical spectra of gauge fields. To appear
- [Yi] Yingst, D.: The Kirchhoff approximation for Maxwell's equations. *Indiana Math. J.* **32** (1983) 543–562
- [Zw1] Zworski, M.: High frequency scattering by a convex obstacle. *Duke Math. J.* (to appear)
- [Zw2] Zworski, M.: Shift of the shadow boundary in high frequency scattering. *Commun. Math. Phys.* (to appear)

# Problems on Limit Sets of Foliations on Complex Projective Spaces

César Camacho

Instituto de Matemática Pura e Aplicada, Estrada Dona Castorina 110  
CEP-22460, Rio de Janeiro, Brasil

We consider differential equations

$$\eta = P(x, y)dy - Q(x, y)dx = 0$$

where  $P(x, y)$  and  $Q(x, y)$  are complex polynomials in the complex variables  $(x, y) \in \mathbf{C}^2$ .

The integrals of this equation are either open Riemann surfaces, passing through points where not both  $P$  and  $Q$  vanish, or singular points  $(x_0, y_0) \in \mathbf{C}^2$  where  $P(x_0, y_0) = Q(x_0, y_0) = 0$ . These integrals define a foliation of  $\mathbf{C}^2$  that extends naturally to a foliation with singularities,  $\mathcal{F}$ , on  $\mathbf{CP}(2)$  the complex projective 2-space. We assume from now on that singularities are isolated.

The first to study these foliations from the local point of view around singularities were C.A.Briot and J.C.Bouquet in 1856 and later by H.Poincaré, P.Painlevé, H.Dulac. Nowadays this local theory is very much developed specially due to the contributions from the Soviet, French and Brazilian schools on this subject, despite the existence of several interesting problems, for instance the developing of a bifurcation theory.

On the other hand, from the global point of view the dynamics of these foliations is far from being understood and it is really in the beginning. Here we wish to pose three problems concerning the most elementary concept of the dynamics of a foliation, that is, its limit set.

As any leaf  $L$  of  $\mathcal{F}$  is open, we define the *limit set* of  $L$  as

$$\lim(L) := \bigcap_{n \geq 1} \overline{L \setminus K_n}$$

where  $K_n \subset K_{n+1} \subset L$  is a sequence of compact subsets of  $L$  such that  $\bigcup_{n \geq 1} K_n = L$ . Then define

$$\lim \mathcal{F} = \overline{\bigcup_L \lim(L)}.$$

So  $\lim \mathcal{F}$  is a closed, invariant subset of  $\mathbf{CP}(2)$ .

**Problem 1.** Classify all foliations whose limit set is analytic.

When  $\lim \mathcal{F}$  is analytic and has dimension zero we have the following

**Theorem** (G. Darboux). *If  $\lim \mathcal{F}$  is finite then  $\mathcal{F}$  admits a meromorphic first integral, i.e. there are polynomials  $f$  and  $g$  such that the leaves of  $\mathcal{F}$  are given by  $f - cg = 0$ ,  $c \in \mathbb{C}$ .*

Notice that by Remmert-Stein theorem  $\lim(L)$  is finite if and only if  $\overline{L}$  is analytic.

A simple example where the limit set has dimension one is the linear foliation given by

$$\mathcal{L} : \lambda x dy - y dx = 0, \quad \lambda \notin \mathbb{R}.$$

Then  $\lim \mathcal{L}$  is a union of three projective lines: the  $x$  and  $y$ -axes, and the line at infinity. The holonomy of each one of these lines is hyperbolic, i.e. it contains an element whose linear part is of the form  $z \mapsto \lambda z$  with  $|\lambda| \neq 1$ . The following theorem shows that foliations with hyperbolic limit set of dimension one are essentially linear.

**Theorem** (C. Camacho, A. Lins N., P. Sad). *Let  $\Lambda = \lim \mathcal{F}$  be an analytic subset of dimension one such that:*

(i) *The holonomy of each irreducible component of  $\Lambda$  is hyperbolic.*

(ii) *The number of separatrices at each singularity is finite.*

*Then there is a rational map  $F$  of  $\mathbb{CP}(2)$  and a linear flow  $\mathcal{L}$  such that  $\mathcal{F} = F^* \mathcal{L}$ .*

*Sketch of Proof.* The proof consists of four parts. The first one contains all the dynamics of the problem.

**I. Lemma.** *Let  $0 \in V \subset \mathbb{C}$  be a neighborhood and  $f, g: V, 0 \rightarrow \mathbb{C}, 0$  holomorphic local diffeomorphisms such that  $|f'(0)| < 1$ . Suppose that for any  $p \in V$  the orbit  $O(p)$  of the pseudogroup generated by  $f$  and  $g$  satisfies  $\overline{O(p)} \setminus O(p) \subset \{0\}$ . Then  $f \circ g = g \circ f$ .*

Thus the holonomy group of each irreducible component of  $\Lambda$  is abelian and linearizable.

**II. The Resolution of  $\mathcal{F}$ .** *Suppose that the singular set of  $\mathcal{F}$  is,  $\text{sing } \mathcal{F} = \{p_1, \dots, p_m\}$ . The resolution of  $\mathcal{F}$  (Theorem of Bendixson-Seidenberg [1, 9]) consists of a proper holomorphic map  $\pi: M \rightarrow \mathbb{CP}(2)$  which is obtained as a certain composition of finitely many quadratic blow up's at the points of  $\text{sing } \mathcal{F}$ , such that if  $D = \pi^{-1} \text{sing } \mathcal{F}$  is the divisor, then:*

(i)  $\pi|_{M \setminus D}$  is a diffeomorphism onto its image.

(ii)  $\tilde{\mathcal{F}} = \pi^* \mathcal{F}$  is a foliation leaving  $D$  invariant with isolated elementary singularities which in local charts  $(x, y)$  have one of the following forms

a) simple:  $\lambda_1 x dy - (\lambda_2 y + \dots) dx = 0$ ,  $\frac{\lambda_1}{\lambda_2} \notin \mathbb{Q}_+$

b) saddle-nodes:  $(x + \dots) dy - y^p dx = 0$ ,  $p \geq 2$ .

**Proposition.** *All components of  $D$  (and so of  $\tilde{\Lambda} = \Lambda \cup D$ ) are hyperbolic and there are no saddle nodes in  $\text{sing } \tilde{\mathcal{F}}$ .*

This follows from the index theorem proved in [2] applied to the resolution of  $\tilde{\mathcal{F}}$ .

**III. Proposition.** *There is a closed meromorphic one form  $\omega$  in  $\mathbf{CP}(2)$  such that:*

- (i) *The polar divisor of  $\omega$ ,  $(\omega)_\infty$ , has order one and is contained in  $\Lambda$ .*
- (ii) *In  $\mathbf{CP}(2) \setminus (\omega)_\infty$ ,  $\omega$  induces  $\mathcal{F}$ .*

The proof of this proposition consists in showing that once all singularities and holonomies of  $\tilde{\Lambda}$  are linearizable it is possible to cover  $\tilde{\Lambda}$  by local chart neighborhoods  $(x_\alpha, y_\alpha) \in U_\alpha$ ,  $\alpha \in A$ , such that if  $\text{sing } \tilde{\mathcal{F}} \cap U_\alpha = \emptyset$  then  $\tilde{\mathcal{F}}|_{U_\alpha}$  is defined by  $dy_\alpha = 0$  and if  $\text{sing } \tilde{\mathcal{F}} \cap U_\alpha \neq \emptyset$  then this intersection is a point  $p \in U_\alpha$  where  $x(p) = y(p) = 0$  and  $\tilde{\mathcal{F}}|_{U_\alpha}$  is given by  $x_\alpha dy_\alpha - \lambda_\alpha y_\alpha dx_\alpha = 0$ . We assume also that if  $U_\alpha \cap U_\beta \neq \emptyset$  then it is simply connected. A form  $\tilde{\omega}$  is then defined in  $U_\alpha U_\alpha$  as  $\tilde{\omega}_\alpha := \frac{dy_\alpha}{y_\alpha}$  in regular neighborhoods  $U_\alpha$  and as  $\tilde{\omega}_\alpha := \frac{dy_\alpha}{y_\alpha} - \lambda_\alpha \frac{dx_\alpha}{x_\alpha}$  in singular neighborhoods  $U_\alpha$ . If  $U_\alpha \cap U_\beta \neq \emptyset$  are two regular neighborhoods it is proven using the existence of a hyperbolic element of the holonomy that the change of coordinates  $y_\beta = L_{\beta\alpha}(y_\alpha)$  is linear, thus  $\tilde{\omega}_\alpha = \tilde{\omega}_\beta$  in  $U_\alpha \cap U_\beta \neq \emptyset$ . Similarly if  $U_\alpha$  is singular,  $U_\beta$  regular and  $U_\alpha \cap U_\beta \neq \emptyset$  one can show that  $\tilde{\omega}_\beta = g_\alpha \cdot \tilde{\omega}_\alpha$  in  $U_\alpha \cap U_\beta$ , where  $g_\alpha$  is holomorphic independent of  $\beta$ .

The form  $\tilde{\omega}$  induces  $\omega = \pi_* \tilde{\omega}$  in a neighborhood of  $\Lambda \subset \mathbf{CP}(2)$ . Now, since  $\mathbf{CP}(2) \setminus \Lambda$  is a Stein manifold this form  $\omega$  extends to  $\mathbf{CP}(2)$  by a generalization of Levi's theorem.

**IV. Construction of  $\mathcal{L}$  and  $F$ .** Here we follow Cerveau-Mattei [4]. The differential form  $\omega$  restricted to  $\mathbf{C}^2$  has the following polar divisor:

$$\Gamma = (\omega)_\infty \cap \mathbf{C}^2 = \bigcup_{j=0}^m \Gamma_j \quad \text{where } \Gamma_j = (f_j = 0)$$

and each  $f_j$  is a polynomial.

Since

$$\mathcal{H}_{DR}^1 \left( \mathbf{C}^2 \setminus \bigcup_{j=0}^m \Gamma_j \right) = \left[ \frac{df_0}{f_0}, \dots, \frac{df_m}{f_m} \right].$$

Then we can show that

$$\omega|_{\mathbf{C}^2} = \sum_{j=0}^m \lambda_j \frac{df_j}{f_j}$$

where

$$\lambda_j = 12\pi i \int_{\gamma_j} \omega$$

and  $\gamma_j$  is a closed path making a simple tour around  $\Gamma_j$  in a small transverse cross section to  $\Gamma_j$ . The holonomy group of  $\Gamma_0$  is generated by  $\{\mu_1^{\ell_1}, \dots, \mu_m^{\ell_m}\} \subset \mathbf{C}^*$  where the  $\ell_j$ 's are integers and  $\mu_j = \exp 2\pi i \lambda_j / \lambda_0$ . Since this group is discrete it is possible to find integers  $u_1, \dots, u_m; v_0, \dots, v_m$  such that for  $\Phi = f_1^{u_1} \dots f_m^{u_m}$  and  $\Psi = f_0^{v_0} \dots f_m^{v_m}$  we have

$$\frac{d\Psi}{\Psi} + \lambda \frac{d\Phi}{\Phi} = \delta \omega|_{\mathbf{C}^2}, \quad \lambda, \delta \in \mathbf{C}^*.$$

Then the map is  $F = (\Phi, \Psi)$  and  $\mathcal{L} : \frac{dy}{y} + \lambda \frac{dx}{x} = 0$ . □

It was shown by Il'iashenko and his students [6, 8] that there are open sets of foliations exhibiting dense leaves, see also Cerveau [5].

The following problem concerns only those foliations which do not admit algebraic leaves.

**Theorem** (A. Lins N. [7]). *For an open and dense set of foliations of  $\mathbf{CP}(2)$  there are no algebraic leaves.*

When no algebraic leaves exist questions about limit sets are more subtle. Along this direction we have.

**Problem 2.** Do nontrivial minimal sets exist?

A minimal set of  $\mathcal{F}$  is a closed, invariant, nonempty subsets of  $\mathbf{CP}(2)$  which is minimal with these three properties and nontrivial means it is not a singularity. Thus the problem above refers to the existence of a leaf of  $\mathcal{F}$  which does not accumulate on  $\text{sing } \mathcal{F}$ . What is known about this problem is the following

**Theorem** (C. Camacho, A. Lins N., P. Sad [3]). *If  $\mathcal{M}$  is a nontrivial minimal set of  $\mathcal{F}$ . Then*

1.  $\mathcal{M}$  is unique.
2. Any leaf accumulates in  $\mathcal{M}$ .
3. Any leaf in  $\mathcal{M}$  has exponential growth, is a hyperbolic Riemann surface and has no parabolic ends.
4. There is no  $\mathcal{F}$ -transverse invariant measure with support in  $\mathcal{M}$ .

Recently Bonatti, Langevin, Moussu, Tischler proved also that there is a leaf in  $\mathcal{M}$  with nontrivial linear holonomy.

Finally we have the following question which intends to establish a relation between these foliations and the iteration of endomorphisms of the Riemann sphere.

**Problem 3** (Sad's Conjecture). Assuming all singularities of  $\mathcal{F}$  generic then  $\lim \mathcal{F}$  is the closure of the separatrices of  $\mathcal{F}$ .

## References

1. Bendixson, I.: Sur les points singuliers d'une équation différentielle linéaire. Ofv. Kongl. Vetenskaps, Akademien Forhandlinger **148** (1895) 81–89
2. Camacho, C., Sad, P.: Invariant varieties through singularities of holomorphic vector fields. Ann. Math. **115** (1982) 579–595
3. Camacho, C., Lins Neto, A., Sad, P.: Minimal sets of foliations in complex projective spaces. Publ. Math. IHES **68** (1988) 187–203
4. Cerveau, D., Mattei, J.F.: Formes intégrables holomorphes singulières. Astérisque **97** (1982)
5. Cerveau, D.: Densité des feuilles de certaines équations de Pfaff à 2 variables. Ann. Inst. Fourier **1** (1983) 185–194
6. Il'iashenko, J.: Global and local aspects of the theory of complex differential equations. Proc. of Int. Cong. of Math. Helsinki 1978, pp. 821–826

7. Lins Neto, A.: Algebraic solutions of polynomial differential equations and foliations in dimension two. *Holomorphic Dynamics Proceedings, Mexico 1986.* (Lecture Notes in Mathematics, vol. 1345.) Springer, Berlin Heidelberg New York 1988
8. Mjuller, B.: On the density of solutions of an equation in  $\mathbb{C}P^2$ . *Mat. Sbornik* **27** (1975) 325–338
9. Seidenberg, A.: Reduction of singularities of the differential equation  $AdY = BdX$ . *Amer. J. Math.* (1968) 248–269



# The Dynamics of Non-uniformly Hyperbolic Systems in Two Variables

*Lennart Carleson*

Department of Mathematics, Royal Institute of Technology, S-10044 Stockholm, Sweden  
and University of California, Los Angeles, CA 90024, USA

One of the most interesting aspects of dynamical systems is to try to understand the nature of ‘chaotic’ behavior and in particular to describe the ‘strange’ attractors. This problem is in a natural way related to the nature of turbulence and there is an overwhelming mass of numerical information but the rigorous mathematical results only apply so far in very special situations.

Nevertheless, there exists today a solid foundation and basic definitions. Let me first recall part of this. We consider a smooth mapping  $x \rightarrow T(x)$  of a compact domain  $\mathcal{D}$  into itself where  $\mathcal{D} \subset \mathbb{R}^n$  or  $\mathbb{T}^n$ . As first realized by Pesin (see e.g. [12]) the starting point should be an invariant measure  $\mu$ , i.e.  $\mu \geq 0$ ,  $\mu(T^{-1}(E)) = \mu(E)$ ,  $\mu(\mathcal{D}) = 1$ . The existence of such a measure is obvious, we can e.g. start from a point mass at a point, push this forward, take an average and a weak limit. The basic fact is now that Liapounov exponents exist a.e. with respect to such a measure. This means that we can split the tangent space  $\oplus E_j(x)$  so that for  $v \in E_j$

$$\lim \frac{\log |DT^n(x)v|}{n} = \lambda_j(x)$$

exists a.e. ( $\mu$ ). In the case of Anosov-mapping this splitting is continuous and  $\lambda_j(x) \neq 0$ , so that expanding and contracting directions are uniformly separated. This may be called the uniformly hyperbolic case. In this case the metric entropy  $h_\mu$  of the pair  $(T, \mu)$  can be written

$$h_\mu = \int \left( \sum \lambda_i(x)^+ \right) d\mu(x)$$

and the conditional measures along the unstable foliation  $W^u(x)$  are absolutely continuous with respect to Lebesgue measure in these directions.

From the point of view of numerical studies a third property of  $(T, \mu)$  is relevant. Let’s call  $\mu$  a physical measure if there exists  $A$  of positive  $n$ -dimensional Lebesgue measure so that for Lebesgue a.a.  $x \in A$  and all  $\varphi \in C_0^\infty(\mathcal{D})$

$$\frac{1}{n} \sum_1^n \varphi(T^j(x)) \rightarrow \int \varphi d\mu.$$

This holds for every ergodic component of  $\mu$  in the Anosov case.

It turns out that all of the above can be generalized to non-uniformly hyperbolic systems. For simplicity assume that  $\lambda_i(x) \neq 0$  a.e. ( $\mu$ ). Briefly stated: the validity of the entropy formula and the absolute continuity of  $\mu$  in unstable directions are equivalent (Ledrappier, Young) [7].  $\mu$  is in this case a physical measure (if it is ergodic) [13]. We call such a  $\mu$  a Sinai-Bowen-Ruelle-measure.

In spite of the generality and completeness, the main question remains open: do (SBR)-measures exist? This question has to be addressed from the beginning and seems to mean that we must build unstable and stable manifolds by hand so that the constructions work Lebesgue almost everywhere. When this has been done one can then read off the Pesin properties of the implied measures.

Let us first discuss the 1-dimensional case. At the previous congress Jakobson [6] gave a report and he also constructed an absolutely continuous invariant measure for the mapping  $T = 1 - a\varphi(x)$  on  $(-1, 1)$  into itself,  $\varphi(t) = t^2$ , for  $a \in$  set of positive Lebesgue measure [5]. The result applies to more general  $\varphi(t)$ , e.g.  $\varphi = |t|^p$ . It is well understood that once we know the increase of the derivative at the critical value

$$D_n = |T^{n'}(1)|$$

then the existence of a (SBR)-measure follows. In fact, Nowicki-van Strien [10], have the following beautiful result: if

$$\sum \varphi^{-1}(D_n^{-1}) < \infty$$

then  $\mu$  exists. (This result should be optimal).

Let us consider somewhat further how one can obtain the increase of  $D_n(a)$ , now considered as a function of  $a$  for, say,  $T = 1 - ax^2$ . Setting  $x_j(a) = T^j(1; a)$  we have  $D_n = \prod_0^{n-1} |-2ax_j(a)|$ . If  $a_0$  corresponds to a situation where 0 is strictly pre-periodic and the resulting cycle is repelling,  $D_n(a_0)$  is obviously exponentially increasing (the simplest case is  $a_0 = 2$ ;  $1 \rightarrow -1 \rightarrow -1 \rightarrow \dots$ ). It is natural to study parameters  $a$  close to  $a_0$  and use the parameter to avoid very small values for  $x_n(a)$ . It is easy to prove that  $|x'_n(a)| \sim D_n(a)$  provided

$$\sum \frac{1}{D_n} < \infty$$

and the larger  $D_n$  the more effectively we can avoid very small values of  $x_n$ . It can be proved that Jakobson's result holds around every  $a_0$  (see [2] for  $a_0 = 2$ ).

The basic property of a dynamical system — as opposed to purely random motions — is the quasi-periodicity. A return of  $x_n$  to  $(-\varepsilon, \varepsilon)$  forces an especially strong repetition of orbits since  $|x_{n+1} - 1| \leq C \cdot \varepsilon^2$  and  $x_{n+j}$  behaves as  $x_j$  as long as  $|x_{n+j} - x_j| < |x_j|/j^2$  (say). Observe also that during this period  $|x_{n+j} - x_j| \sim D_j x_n^2$ . If we slightly weaken the sufficient condition  $D_j x_n^2 < 1/j^2$  to  $< 1$  we obtain the following model (see [4]) for the amount of randomness in the sequence  $D_n(a)$ .

Let  $\{\Delta_n(\omega)\}_0^\infty$  be the following stochastic process.

- (a)  $\Delta_0(\omega) \equiv 1$ ;
- (b) there exist stopping times  $\{n_j(\omega)\}_1^\infty$ ,  $n_1(\omega) \equiv 1$ ;

(c) If  $n$  is a stopping time, choose a number  $t_n$  at random with uniform distribution in  $(0, A)$  and independent of the past and define

$$\Delta_{n+j}(\omega) = \Delta_{n-1}(\omega) t_j \Delta_j(\omega)$$

for  $j = 0, 1, \dots, m-1$  as long as

$$t^2 \Delta_j(\omega) \leq 1.$$

If  $m < \infty$  the next stopping time is  $n+m$ .

It is an easy fact that if  $A > e$  then

$$\lim_{n \rightarrow \infty} \frac{\log \Delta_n}{n} > 0$$

with probability 1.

Interpreting parameter space as probability space and  $\Delta_n = D_n$  this statement contains all the essentials of the chaotic behavior for these one-dimensional systems.

Let us now turn to the two-variable case. J. Palis [11] has proposed the following scenario for the creation of strange attractors (i.e. an invariant set  $\Lambda$ ,  $W^s(\Lambda) \supset$  a neighborhood of  $\Lambda$ , with a dense orbit with positive Liapounov exponent). Let  $T(x; a)$  depend on a parameter and let 0 be a fixed point which is hyperbolic (one eigenvalue  $> 1$ , one  $< 1$ , the product  $< 1$ ). Let  $W^s(a)$ ,  $W^u(a)$  be the stable and unstable manifolds through 0 and suppose  $W^s(a_0)$  is tangent to  $W^u(a_0)$  at some point  $p_0$ . Then, for  $a$  close to  $a_0$ ,  $T(x; a)$  has, for a set of positive measure of the parameter, a strange attractor. This has in fact been proved recently by Mora-Viana [8], see Fig. 1.

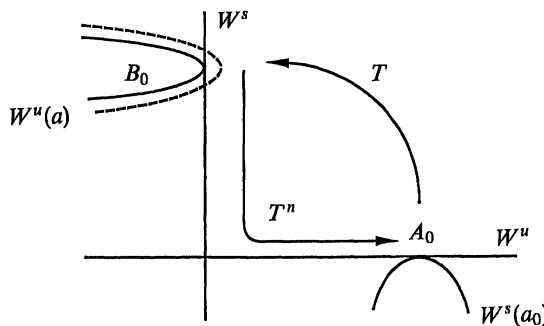


Fig. 1

It turns out that the basic properties of  $T$  show up already for the special (Hénon) map

$$T_0 : \begin{cases} 1 - ax^2 + y & 0 < b < b_0 \\ bx & a_0 < a < 2 \end{cases}$$

with  $b_0$  arbitrarily small and  $a_0$  close to 2. The reason is that some suitable high iterate  $T^n$  of  $T$  maps a small box close to  $p_0$  into itself and the basic behavior of  $T^n$  is obtained from its second order terms and the eigenvalues at 0. Let us therefore try to understand the special case  $T_0$ , see Fig. 2.

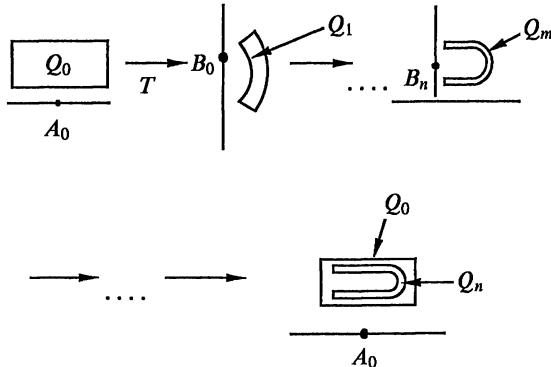


Fig. 2

$T_0(x, y)$  has a fixed point  $F$  at  $\sim (\frac{1}{2}, \frac{b}{2})$  and the eigenvalues are  $\sim (-2, \frac{b}{2})$ . The unstable manifold  $W^u$  looks like Figure 3 (see next page), and it is easy to see that if an attractor exists it is  $\subset \overline{W^u}$ . We would expect  $\overline{W^u}$  to be the strange attractor. Since  $b$  is small we expect  $\|DT_0^n\| \sim D_n$  corresponding to the case  $b = 0$ . Most directions would then expand for  $T_0^n$  but since  $\det(T_0^n) = (-b)^n$  some very strong contracting direction must exist. We need to find this and to make a foliation along these directions. The expanding foliation is clearly  $\overline{W^u}$  and we need most leaves to be long. The 'critical' values are those  $c \in W^u$  for which the tangent direction of  $W^u$  = the contracting direction and this set is located in  $(-b, b)$  as indicated below (Only the intersection with  $\overline{W^u}$  is relevant and this is formally a Cantor set). It now turns out that existence of contracting directions depends on the simultaneous existence of expanding directions but the accuracy is increased by powers in  $b$  and the situation therefore allows for a boot-strap-argument. We prove in the end [3] that all critical points are uniformly expanding after parameter exclusion. This amounts to using the 1-dimensional argument and to verifying, for the Cantor set of possible starting points  $c$ ,

$$\|DT_0^n(c)\| \geq e^{yn}.$$

We now have the setup for a foliation and for the understanding of  $W^u$ . A (SBR)-measure can be constructed [1] starting from arc length on a sheet of  $W^u$  and pushing this forward. Most mass will be located on rather long arcs  $\subset W^u$ .

Let me elaborate some more on this constructive aspect of the Pesin theory.

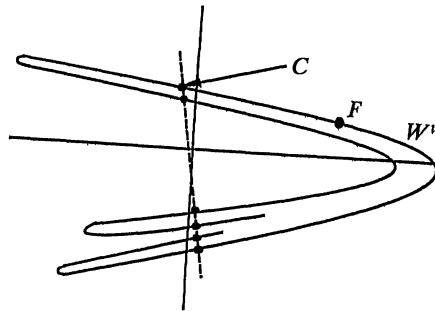


Fig. 3

$$DT_0 = \begin{pmatrix} -2ax_0 & 1 \\ b & 0 \end{pmatrix}$$

transfers a vector  $(1, q)$  of slope  $q$  at  $(x_0, y_0)$  to a vector of slope  $q_1$  where

$$q = 2ax_0 + \frac{b}{q_1}$$

and

$$\|DT(1, q)\|^2 = (-2ax_0 + q)^2 + b^2.$$

For  $q = q_0 = 2ax_0$  we have length  $b$  so this direction is almost maximally contracting. If we choose

$$q_1 = 2ax_1$$

we again contract and get in the general case

$$q = 2ax_0 + \frac{b}{\frac{2ax_1 + \frac{b}{q_2}}{\ddots + 2ax_{n-1}}}.$$

It is natural to define

$$\bar{q}_n = 2ax_0 + \frac{b}{2ax_1 + \frac{b}{\ddots + 2ax_{n-1}}}$$

and call this the contracting direction and

$$V_n : \quad y' = \bar{q}_n(x, y)$$

the  $n$ -th order contracting vector field. The  $\bar{q}_n$  are very unstable but if we have expansion along  $(x_0, \dots, x_{n-1})$

$$|\bar{q}_n(x, y) - \bar{q}_{n-1}(x, y)| \leq b^{n-1}.$$

We now obtain a stable foliation by integrating along the vector field  $V_n$  where we get contraction by  $b^n$ . The expansion is preserved under this small perturbation and we obtain in the limit  $W^s(x_0, y_0)$  for Lebesgue a.a. points  $(x_0, y_0)$  along  $W^u$ . Whether this holds for Lebesgue  $(dx dy)$  a.e. point in  $\mathcal{D}$  remains open but is probably true.

Let me summarize what we now know in a theorem.

*Given  $b > 0$ , sufficiently small, there exists a set  $E$  of positive measure so that for  $a \in E$*

- (i)  $T_0(x, y)$  has a strange attractor
- (ii)  $T_0(x, y)$  admits a unique (SBR)-measure
- (iii) for any  $a \in E$ , there exists  $a' \notin E$ , so that any  $\varepsilon > 0$   $|a - a'| < \varepsilon$  and  $T(x, y, a')$  has an attractive cycle.

It is still not known if for  $a \in E$  almost all (Lebesgue) points in  $\mathcal{D}$  generate the (SBR)-measure but this seems very likely.

Let me end by some comments.

The original Hénon simulation concerned the parameter values  $(1.4, 0.3)$ . It is most likely that the Mora-Viana-result applies and we have strange attractors arbitrarily nearby. However, the most natural conjecture here is that the problem of proving existence of a strange attractor for a particular parameter value is in some rigorous sense undecidable. To develop this aspect of (even one-dimensional) dynamics seems an interesting task.

The approach can be generalized to  $n$  dimensions when  $(n - 1)$  eigenvalues are very small and one eigenvalue is  $> 1$ ; so that the unstable manifold is 1-dimensional [9]. The case of  $> 1$ -dimensional  $W^u$  has not been studied. One can see possibilities of progress also then.

It was perhaps not so clear from the presentation but it is of central importance for the argument that the Cantor set of critical points is very thin. This made it possible to consider every critical point independently and obtain uniform expansion. We could then disregard the distribution of the orbit in the "vertical" direction. However, for an interesting case such as the standard map on  $\mathbb{T}^2$ :

$$\begin{cases} x' = k \sin x - 2x - y \\ y' = x \end{cases}$$

the critical set basically is two 'curves' ( $x \sim \frac{\pi}{2}, \frac{3\pi}{2}$ ) and the distribution along these 'curves' must be understood. The problem here remains open even if  $k$  is very large.

## References

1. Benedicks, M., Carleson, L.: On iterations of  $1 - ax^2$  on  $(-1, 1)$ . Ann. Math. **122** (1985) 1–25
2. Benedicks, M., Carleson, L.: The dynamics of the Hénon map. Ann. Math. **132** (1990) 629–725
3. Benedicks, M., Young, L.-S.: To appear.

4. Carleson, L.: A model for chaotic dynamics. To appear in Lecture Notes in Mathematics. Springer, Berlin Heidelberg New York
5. Jakobson, M.: Absolutely continuous invariant measures for one-parameter families of one-dimensional maps. Comm. Math. Phys. **81** (1981) 39–88
6. Jakobson, M.: Families of one-dimensional maps and nearby diffeomorphisms. Proc. Int. Congr. ICM-86 Berkeley 1986, vols. 1, 2. Amer. Math. Soc., pp. 1150–1160
7. Ledrappier, F., Young, L.-S.: The metric entropy of diffeomorphisms, I, II. Ann. Math. **122** (1985) 509–574
8. Mora, L., Viana, M.: Abundance of strange attractors. To appear in Acta Math.
9. Mora, L., Viana, M.: To appear
10. Nowicki, J., van Strien, S.: Invariant measures exist under a summability condition for unimodal maps. Preprint Delft
11. Palis, J., Takens, F.: Homoclinic bifurcations. Cambridge Univ. Press. To appear
12. Pesin, Ya. B.: Characteristic Liapunov exponents and smooth ergodic theory. Russ. Math. Surv. **32** (4) (1977) 55–114
13. Pugh, C., Shub, M.: Ergodic attractors. Trans. Amer. Math. Soc. **312** (1989) 1–54



# The Acceleration Operators and Their Applications to Differential Equations, Quasianalytic Functions, and the Constructive Proof of Dulac's Conjecture

*Jean P. Ecalle*

Mathématiques, Bâtiment 425, Université de Paris-Sud, Centre d'Orsay  
F-91405 Orsay, France

We introduce (§§1,2) a general apparatus (resurgence, alien derivations, acceleration, etc.) that enables one to study and resum most divergent expansions of natural origin. We then proceed to give three select applications (§§3–5).

## §1. Resurgent Functions, Alien Derivations and Medianization

### The Singularity Algebras $\mathcal{S}(S_\theta)$ and $\mathcal{S}^{\text{int}}(S_\theta)$

Let  $\mathbb{C}$  be the Riemann surface of  $\log \zeta$  and  $S_\theta$  (resp.  $S_{\theta_1, \theta_2}$ ) the semi-axis  $\arg \zeta = \theta$  (resp. the sector  $\theta_1 \leq \arg \zeta \leq \theta_2$ ). A *major*  $\check{\phi}$  on  $S_\theta$  is an analytic germ defined on  $S_{\theta-2\pi, \theta}$  close to 0 (“at the root of  $S_{\theta-2\pi, \theta}$ ”). The *minor*  $\hat{\phi}$  of  $\check{\phi}$  is defined at the root of  $S_\theta$  by  $\hat{\phi}(\zeta) = \check{\phi}(\zeta) - \check{\phi}(\zeta \cdot e^{-2\pi i})$ . A *singularity* of direction  $\theta$  is a class  $\check{\phi}$  of majors  $\check{\phi}$  modulo the space of regular (i.e. holomorphic) functions at 0. A singularity  $\check{\phi}$  is said to be *integrable* iff:

$$(1.1) \quad \begin{cases} \zeta \check{\phi}(\zeta) \rightarrow 0 \text{ as } \zeta \rightarrow 0 & \text{on } S_{\theta-2\pi, \theta} \text{ and} \\ \int_0^{\zeta_0} |\hat{\phi}(\zeta)| |d\zeta| < +\infty & \text{for } \zeta_0 \in S_\theta \text{ close to 0.} \end{cases}$$

Integrable singularities are fully determined by their minor.

For any two classes  $\check{\phi}_1, \check{\phi}_2 \in \mathcal{S}(S_\theta)$  and  $u \in S_{\theta-2\pi}$  close to 0, the class  $\check{\phi}_3$  of the *major*  $\check{\phi}_{3,u}$  defined by

$$(1.2) \quad \begin{aligned} \check{\phi}_{3,u}(\zeta) &= \int_u^{\zeta-u} \check{\phi}_1(\zeta_1) \check{\phi}_2(\zeta - \zeta_1) d\zeta, \\ (u \in S_{\theta-2\pi}; \zeta &\text{ and } \zeta - u \in S_{\theta-2\pi, \theta-\pi}) \end{aligned}$$

depends neither on  $u$  nor on the choice of  $\check{\phi}_i$  in  $\check{\phi}_i$ . The convolution  $\check{\phi}_1 * \check{\phi}_2 = \check{\phi}_3$  thus defined turns  $\mathcal{S}(S_\theta)$  into a commutative algebra. The space  $\mathcal{S}^{\text{int}}(S_\theta)$  of integrable singularities is a subalgebra and its convolution reduces to:

$$(1.3) \quad \hat{\phi}_1 * \hat{\phi}_2(\zeta) = \int_0^\zeta \hat{\phi}_1(\zeta_1) \hat{\phi}_2(\zeta - \zeta_1) d\zeta_1 \quad (\zeta, \zeta_1, \zeta - \zeta_1 \text{ on } S_\theta \text{ and close to 0}).$$

## The Resurgence Algebras $\mathcal{R}(S_\theta)$ and $\mathcal{R}^{\text{int}}(S_\theta)$

The subspace  $\mathcal{R}(S_\theta) \subset \mathcal{S}(S_\theta)$  of all  $\check{\phi}$  whose minor  $\hat{\phi}$  can be analytically continued along any path that keeps close to  $S_\theta$  (without going back) and bypasses to the right or to the left all intervening singular points  $\omega_i \in S_\theta$ , is closed under convolution. For any sequence  $\varepsilon_i$  of  $\pm$  signs and any  $\zeta$  in  $\] \omega_r, \omega_{r+1} [$ , denote by  $\varphi_{\omega_1, \dots, \omega_r}^{\varepsilon_1, \dots, \varepsilon_r}(\zeta)$  the determination of  $\hat{\phi}(\zeta)$  obtained by starting from 0 and bypassing each  $\omega_i$  to the right (resp. left) if  $\varepsilon_i = +$  (resp.  $-$ ). The space of all  $\check{\phi}$  in  $\mathcal{R}(S_\theta)$  whose minors  $\hat{\phi}$  have all their determinations  $\hat{\phi}_{\omega_i}^{\varepsilon_i}$  integrable on their segment of definition  $\] \omega_r, \omega_{r+1} [$ , constitutes a subalgebra  $\mathcal{R}^{\text{int}}(S_\theta)$ . We call  $\mathcal{R}(S_\theta)$  (resp.  $\mathcal{R}^{\text{int}}(S_\theta)$ ) the algebra of resurgent functions (resp. integrable resurgent functions) of direction  $\theta$ .

### Alien Derivations and Medianization

For any finite sequence  $\varepsilon_i = \pm$  denote by  $p$  (resp.  $q$ ) the number of  $+$  (resp.  $-$ ) signs and consider the weights:

$$(1.4) \quad \delta^{\varepsilon_1, \dots, \varepsilon_{r-1}} = \delta_{p,q} = \frac{p! q!}{(p+q+1)!}; \quad \lambda^{\varepsilon_1, \dots, \varepsilon_r} = \lambda_{p,q} = \frac{(2p)! (2q)!}{4^{p+q} p! q! (p+q)!}$$

For any  $\omega \in S_\theta$  the operator  $\Lambda_\omega$  of  $\mathcal{R}(S_\theta)$  onto itself defined by:

$$(1.5) \quad \Lambda_\omega : \check{\phi} \longmapsto \check{\phi}_\omega \quad \text{with} \quad \check{\phi}_\omega(\zeta) = \sum_{\varepsilon_i} \delta^{\varepsilon_1, \dots, \varepsilon_{r-1}} \hat{\phi}_{\omega_1, \dots, \omega_{r-1}, \omega}^{\varepsilon_1, \dots, \varepsilon_{r-1}, +}(\zeta + \omega)$$

(for  $\zeta$  on  $S_\theta$  and close to 0) is a *derivation* of the algebra  $\mathcal{R}(S_\theta)$ :

$$(1.6) \quad \Lambda_\omega(\check{\phi}_1 * \check{\phi}_2) = (\Lambda_\omega \check{\phi}_1) * \check{\phi}_2 + \check{\phi}_1 * (\Lambda_\omega \check{\phi}_2).$$

We call  $\Lambda_\omega$  the *alien derivation* with index  $\omega$ . On  $\mathcal{R}^{\text{int}}(S_\theta)$  it reduces to:

$$(1.7) \quad \begin{aligned} \Lambda_\omega : \hat{\phi} &\longmapsto \hat{\phi}_\omega \\ \text{with } \hat{\phi}_\omega(\zeta) &= \sum_{\varepsilon_i} \delta^{\varepsilon_1, \dots, \varepsilon_{r-1}} \left\{ \hat{\phi}_{\omega_1, \dots, \omega_{r-1}, \omega}^{\varepsilon_1, \dots, \varepsilon_{r-1}, +}(\zeta + \omega) - \hat{\phi}_{\omega_1, \dots, \omega_{r-1}, \omega}^{\varepsilon_1, \dots, \varepsilon_{r-1}, -}(\zeta + \omega) \right\}. \end{aligned}$$

Along with the natural derivation  $\hat{\partial} : \check{\phi} \mapsto -\zeta \check{\phi}$ , the  $\Lambda_\omega$  finitely generate all continuous derivations of  $\mathcal{R}(S_\theta)$ .

Similarly, the operator med (“medianization”):

$$(1.8) \quad \text{med} : \hat{\phi}(\zeta) \longmapsto \text{med } \hat{\phi}(\zeta) = \sum_{\varepsilon_i} \lambda_{\omega_1, \dots, \omega_r}^{\varepsilon_1, \dots, \varepsilon_r}(\zeta) \quad (\text{if } \omega_r < \zeta < \omega_{r+1})$$

is a homomorphism of the algebra  $\mathcal{R}^{\text{int}}(S_\theta)$  into the algebra  $\mathbb{L}^{\text{int}}(S_\theta)$  of univalued, locally integrable functions on  $S_\theta$ :

$$(1.9) \quad \text{med } (\hat{\phi}_1 * \hat{\phi}_2) = (\text{med } \hat{\phi}_1) * (\text{med } \hat{\phi}_2) \quad \text{with } * \text{ as in (1.3).}$$

Medianization has the added advantage of preserving realness: if the *germ*  $\hat{\phi}$  is real, so is the *function*  $\text{med } \hat{\phi}$ .

### The Resurgence Algebras $\mathcal{R}$ and $\mathcal{R}^{\text{int}}$

The convolution algebra of all classes  $\check{\phi}$  which, along with their successive alien derivatives  $\Delta_{\omega_1} \dots \Delta_{\omega_l} \check{\phi}$  ( $\forall \omega_i \in \mathbb{C}$ ) belong to all  $\mathcal{R}(S_\theta)$  (resp. all  $\mathcal{R}^{\text{int}}(S_\theta)$ ) is known as the general algebra  $\mathcal{R}$  (resp.  $\mathcal{R}^{\text{int}}$ ) of resurgent functions (resp. integrable r.f.). Their *majors*  $\check{\phi}$  are defined in spiral-like neighbourhoods of  $0$  on  $\mathbb{C}$  and their *minors*  $\hat{\phi}$  can be continued (starting from  $0$  and bypassing intervening singular points) along any split line on  $\mathbb{C}$ .

Resurgent functions of natural origin tend to reproduce themselves at their singular points. This self-reproduction is exactly described by *resurgence equations* linking  $\check{\phi}$  to its *alien derivatives*  $\Delta_{\omega} \check{\phi}$ . By a slight abuse, we often extend the label of “resurgent function” to those power series whose Borel transforms (see §2) belong to  $\mathcal{R}$ . See [1–3].

## §2. The Acceleration Operators

### The Laplace Transform $\mathcal{L}$ and the Borel Transforms $\mathcal{B}$ and $\tilde{\mathcal{B}}$

$$(2.1) \quad \mathcal{L} : \hat{\phi}(\zeta) \longmapsto \varphi(z) = \int_0^{+\infty} e^{-\zeta z} \hat{\phi}(\zeta) d\zeta$$

$$(2.2) \quad \mathcal{B} : \varphi(z) \longmapsto \hat{\phi}(\zeta) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} e^{z\zeta} \varphi(z) dz$$

The classical *Laplace transform*  $\mathcal{L}$  is a homomorphism of the convolution algebra  $\mathbb{L}_{\text{exp}}^{\text{int}}(\mathbb{R}^+)$  of all univalued, locally integrable functions on  $\mathbb{R}^+$  with (at most) exponential growth at  $+\infty$ , into the multiplicative algebra  $\mathbb{B}$  of holomorphic germs bounded in half planes  $\text{Re } z \geq x_0$ . Its inverse  $\mathcal{B}$  is known as the Borel transform. For each formal series  $\tilde{\varphi}(z) = \sum \varepsilon_n(z)$  whose general term  $\varepsilon_n \in \mathbb{B}$  has a Borel transform  $\hat{\varepsilon}_n$ , we have a notion of *formal* (or term-wise) Borel transform:

$$(2.3) \quad \tilde{\mathcal{B}} : \tilde{\varphi}(z) = \sum \varepsilon_n(z) \longmapsto \hat{\phi}(\zeta) = \sum \hat{\varepsilon}_n(\zeta) \quad \left( \text{e.g. } \sum a_n z^{-n} \longmapsto \sum \frac{a_n}{\Gamma(n)} \zeta^{n-1} \right).$$

### The Acceleration Operators and Their Kernels

An acceleratrix is a function  $F$  holomorphic in a neighbourhood of  $\infty \in \mathbb{C}$ , real positive on  $\mathbb{R}^+ = S_0$  and such that for  $z \rightarrow \infty$ :

$$(2.4) \quad x^{-1}F(x) \rightarrow 0; \quad \delta F(z) \sim \delta F(x); \quad \delta^2 F(z) \sim \delta^2 F(x)$$

with  $0 < x \rightarrow \infty$ ;  $z = xe^{i\theta}$  ( $\theta$  fixed in  $\mathbb{R}$ ) and:

$$(2.4\text{bis}) \quad \delta F(z) = zF'(z)F(z); \quad \delta^2 \varphi(z) = 1 + zF''(z)/F'(z) - zF'(z)/F(z).$$

The co-acceleratrix  $G$  of  $F$  is the germ  $G$  defined on  $[+0, \dots]$  by:

$$(2.5) \quad G(f(z)) \equiv F(z) - zf(z) \text{ with } f(z) = F'(z) \quad (G(\zeta) \rightarrow +\infty \text{ as } \zeta \rightarrow 0).$$

Borel-Laplace takes the multiplicative endomorphism  $\mathcal{C}_F : \varphi_1(z_1) \mapsto \varphi_2(z_2) \equiv \varphi_1(F(z_2))$  of  $\mathbb{B}$  into a convolution endomorphism  $\widehat{\mathcal{C}}_F$  of  $\mathbb{L}_{\exp}^{\text{int}}(\mathbb{R}^+)$ :

$$(2.6) \quad \widehat{\mathcal{C}}_F = \mathcal{B.C}_F.\mathcal{L} : \hat{\varphi}_1(\zeta_1) \longmapsto \hat{\varphi}_2(\zeta_2) = \int_0^{+\infty} C_F(\zeta_2, \zeta_1) \hat{\varphi}_1(\zeta_1) d\zeta_1.$$

$\widehat{\mathcal{C}}_F$  is known as the *acceleration*  $z_1 \rightarrow z_2$ . Its integral kernel is given by:

$$(2.7) \quad C_F(\zeta_2, \zeta_1) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \exp(\zeta_2 z_2 - \zeta_1 F(z_2)) dz_2$$

and it has faster-than-exponential decrease in  $\zeta_1$ :

$$(2.8) \quad \log C_F(\zeta_2, \zeta_1) \sim -\zeta_1 G(\zeta_2/\zeta_1) \quad \text{for } 0 < \zeta_2 \text{ fixed and } \zeta_1 \rightarrow +\infty.$$

Therefore, the natural domain of definition of  $\widehat{\mathcal{C}}_F$  is much larger than  $\mathbb{L}_{\exp}^{\text{int}}(\mathbb{R}^+)$ ; it is the convolution algebra  $\mathbb{L}_{F\text{-acc}}^{\text{int}}(\mathbb{R}^+)$  of all functions  $\hat{\varphi}_1(\zeta_1)$  with  $F$ -acceleratable growth, i.e. for which there exists  $c > 0$  such that:

$$(2.9) \quad |\hat{\varphi}_1(\zeta_1)| \leq \text{Cst.}/C_F(c, \zeta_1) \quad \text{or} \quad \log |\hat{\varphi}_1(\zeta_1)| \leq \zeta_1 G(c/\zeta_1) \text{ as } \zeta_1 \rightarrow +\infty.$$

The largest such  $c$  is the *acceleration abscissa* of  $\hat{\varphi}_1$ .

### Strong, Moderate, Weak Accelerations

*Strong accelerations* ( $\log z_2 / \log z_1 \rightarrow +\infty$ ) have kernels with *slightly* over-exponential decrease in  $\zeta_1$  and yield germs  $\hat{\varphi}_2(\zeta_2)$  which are defined in a spiral-like neighbourhood of  $0 \in \mathbb{C}$  with infinite aperture. *Moderate accelerations* ( $\log z_2 / \log z_1 \rightarrow 1/\alpha$  with  $0 < \alpha < 1$ ) have kernels decreasing like  $\exp(-c_\alpha \zeta_1^{1/\beta} \zeta_2^{-\alpha/\beta})$  with  $\beta = 1 - \alpha$  and they yield germs  $\hat{\varphi}_2(\zeta_2)$  which are defined in a sector of aperture  $\pi\beta/\alpha$ . *Weak accelerations* ( $\log z_2 / \log z_1 \rightarrow 1$  but  $z_2/z_1 \rightarrow 1$ ) have very fast decreasing kernels but they yield germs  $\hat{\varphi}_2(\zeta_2)$  which are defined only at the root of  $\mathbb{R}^+ = S_0$  and are usually non-analytic, but only Denjoy-quasianalytic (cf. §4). Of great practical importance are the elementary accelerations  $z_1 \rightarrow z_2 = \exp(\sigma z_1)$  and  $z_1 \rightarrow z_2 = z_1^{1/\alpha}$  with their respective kernels:

$$(2.10) \quad C_F(\zeta_2, \zeta_1) = (\zeta_2)^{\zeta_1/\sigma-1} / \Gamma(\zeta_1/\sigma)$$

$$(2.11) \quad C_F(\zeta_2, \zeta_1) = \zeta_2^{-1} C_\alpha(X)$$

with  $X = \zeta_1^{1/\beta} \zeta_2^{-\alpha/\beta} \quad (0 < \alpha, \beta < 1, \alpha + \beta = 1)$

$$(2.11\text{bis}) \quad C_\alpha(X) \sim (c/2\pi)^{1/2} \cdot X^{1/2} \cdot \exp(-cX)$$

when  $X \rightarrow \infty$  with  $\operatorname{Re} X > 0$  ( $c = \alpha^{\alpha/\beta} \beta$ ).

### Accelero-summability

A formal series  $\tilde{\varphi}(z) = \sum e_n(z)$  is said to be accelero-summable with sum  $\varphi(z)$  and critical times  $z_1, z_2, \dots, z_r$  if it can be subjected to the following operations (algebra homomorphisms):

$$\begin{array}{ccc} \tilde{\varphi}_1(z_1) = \tilde{\varphi}(z) & \xrightarrow{\quad\quad\quad} & \varphi(z) = \varphi_r(z_r) \\ \downarrow \tilde{\mathcal{B}} & & \uparrow \mathcal{L} \\ \hat{\varphi}_1(\zeta_1) \xrightarrow{12} \hat{\varphi}_2(\zeta_2) \xrightarrow{23} \hat{\varphi}_3(\zeta_3) \rightarrow \cdots \rightarrow \hat{\varphi}_r(\zeta_r) & & \text{convolution algebras} \end{array}$$

with arrows  $(i, i+1)$  denoting the acceleration  $z_i \rightarrow z_{i+1}$ . See [4, 5, 7].

### §3. Acceleration Applied to Many-Levelled Differential Systems

Resurgent functions are truly ubiquitous. They arise as formal solutions of differential equations (or difference equations, or general functional equations) with analytic coefficients, or of systems of such equations. They occur in the study (normalization, conjugacy, iteration) of local analytic objects, chiefly: local singular vector fields and local diffeomorphisms. Again, most expansions in a “singular parameter” (such as the Planck constant in the Schrödinger equation) turn out to be divergent and resurgent. Indeed, it is no exaggeration to claim that most divergent expansions met with in actual life are not only resurgent but also summable by  $\mathcal{LB}$  (one critical time) or by  $\mathcal{LC}_{F_{r-1}} \dots \mathcal{C}_{F_r} \mathcal{B}$  ( $r$  critical times;  $r \geq 2$ ). The former case ( $r = 1$ ) is by far the more common. Criticity  $r \geq 2$  occurs only in connection with objects of a certain complexity. Thus it never arises with vector fields (resp. diffeomorphisms) in less than 3 (resp. 2) dimensions.

We shall describe that phenomenon (namely,  $r \geq 2$ ) in the case of a *many-levelled* but formally separable differential system, because it illustrates all the relevant analysis while keeping formal complications down to a minimum. So, consider a local analytic system (3.1) that is formally conjugate to the normal system (3.2) under transformation (3.3).

$$(3.1) \quad \frac{1}{p_i} t^{1+p_i} \dot{x}_i + \lambda_i x_i = b_i(t, x_1, \dots, x_v) \in \mathbb{C}\{t, x_1, \dots, x_v\} \quad (i = 1, \dots, v)$$

$$(3.2) \quad \frac{1}{p_i} t^{1+p_i} \dot{y}_i + \lambda_i y_i = 0 \quad \left( i = 1, \dots, v ; \lambda_i \in \mathbb{C}^* ; p_i \in \mathbb{N}^* \right)$$

$$(3.3) \quad x_i = h_i(t, y_1, \dots, y_v) \in \mathbb{C}[[t, x_1, \dots, x_v]] \quad (i = 1, \dots, v).$$

Let  $q_1 < q_2 < \dots < q_r$  be the distinct values taken by the *levels*  $p_i$  and assume for simplicity's sake that the various  $\lambda_i$  attached to any given level display neither resonance nor quasiresonance (i.e. the combinations  $\sum n_i \lambda_i$  neither vanish nor do

they get abnormally close to 0). Under those mild genericity assumptions, the formal integral  $x(t, u)$  of system (3.1):

$$(3.4) \quad x(t, u) = \sum u^n E^n(t) \varphi^n(t) \\ \left( n \in \mathbb{N}^v ; u^n = \prod u_i^{n_i} ; E^n(t) = \exp \sum n_i \lambda_i t^{-p_i} ; \varphi^n \in (\mathbb{C}[[t]])^v \right)$$

obtained by plugging into (3.3) the elementary solution  $y_i = u_i \exp(\lambda_i t^{-p_i})$  of (3.2), can be shown to be convergent in  $u$  and divergent in  $t$ , but resurgent and accelero-summable with critical times  $z_i = t^{-q_i}$  ( $i = 1, \dots, r$ ). This compact statement translates into the following. Let  $\theta$  be a multipolarization, i.e. a choice of angles  $\theta_1, \dots, \theta_r$  satisfying the self-compatibility condition:

$$(3.5) \quad \left| \frac{\theta_i}{q_i} - \frac{\theta_{i+1}}{q_{i+1}} \right| < \frac{\pi}{2} \left( \frac{1}{q_i} - \frac{1}{q_{i+1}} \right) \quad (1 \leq i \leq r-1 ; \theta_i \in \mathbb{R}).$$

Further, let  $\Omega_i$  be the set of all  $\omega \in \mathbb{C}$  whose projection  $\hat{\omega}$  on  $\mathbb{C}$  is of the form  $\sum n_j \lambda_j$  with  $p_j = q_i$  and  $n_j \in \mathbb{N}$  (or  $n_j = -1$  for one  $j$  at most). Then each component  $\varphi^n(t) = \varphi_1^n(z_1)$  of  $x(t, u)$  has a Borel transform  $\hat{\varphi}_1(\zeta_1)$  with only isolated singularities and a growth rate not exceeding  $\exp(\text{cst.} |\zeta_1|^{q_2/q_2 - q_1})$ . Thus it has accelerable growth for the acceleration  $z_1 \rightarrow z_2$  taken along any axis  $\arg \zeta_1 = \theta_1$  that avoids singularities. The corresponding accelerate  $\hat{\varphi}_2(\zeta_2 \parallel \theta)$  can be analytically continued within a sector  $\hat{S}_2$  containing (at least) all directions  $\theta_2$  linked to  $\theta_1$  by (3.5). In that sector it possesses only isolated singularities and has accelerable growth for the acceleration  $z_2 \rightarrow z_3$ . Thus we get a succession of accelerates  $\hat{\varphi}_i(\zeta_i \parallel \theta)$ , the last of which (for  $i = r$ ) has exponential growth and can be laplaced along any semi-axis  $\theta_r$  compatible with  $\theta_{r-1}$ , yielding the sought-after sum  $\varphi_r(z_r \parallel \theta) = \varphi(t \parallel \theta)$ .

Moreover, each  $i$ -th Borel transform  $\hat{x}(\zeta_i, u \parallel \theta)$  satisfies the so-called Bridge Equation, which reads:

$$(3.6) \quad \dot{\Delta}_\omega \hat{x}(\zeta_i, u \parallel \theta) = \mathbf{A}_{\omega \parallel q_i, \theta} \cdot \hat{x}(\zeta_i, u \parallel \theta) \quad \left( \dot{\Delta}_\omega = e^{-\omega z} \Delta_\omega ; \omega \in \Omega_i \cap \hat{S}_i \right)$$

$$(3.6\text{bis}) \quad \mathbf{A}_{\omega \parallel q_i, \theta} = u^{n(\omega)} \left\{ \sum_{p_j \geq q_i} \mathbf{A}_{\omega \parallel q_i, \theta}^j(u) \cdot u_j \frac{\partial}{\partial u_j} + \sum_{p_j < q_i} \mathbf{A}_{\omega \parallel q_i, \theta}^j(u) \cdot \frac{\partial}{\partial u_j} \right\}$$

$$(3.6\text{ter}) \quad \left\{ \begin{array}{l} \dot{\omega} = \sum_{p_j = q_i} n_j(\omega) \lambda_j ; \quad u^{n(\omega)} = \prod u_j^{n_j(\omega)} ; \\ \mathbf{A}_{\omega \parallel q_i, \theta}^j \in \mathbb{C}[[u_k \text{ with } p_k < q_i]] \end{array} \right.$$

and which describes in compact form the resurgence properties of all the  $\hat{\varphi}_i(\zeta_i \parallel \theta)$ . The Bridge Equation holds, in some form or other, for all local objects. It says in effect that alien derivations  $\dot{\Delta}_\omega$  act on formal integrals like ordinary differential operators  $\mathbf{A}_\omega$ , while at the same time enabling one to calculate those  $\mathbf{A}_\omega$ .

Here, the components of the  $\mathbf{A}_{\omega \parallel q_i, \theta}$  relative to each level  $q_i$  are formal power series in all the parameters  $u_k$  attached to the lower levels  $p_k < q_i$ . For  $i = 1$ , they are scalar-valued. Lastly, and crucially, these differential operators, taken together, constitute a complete set of *analytic invariants* of the system (3.1).

The proof [6, 8] relies heavily on the study of the operators:

$$(3.7) \quad D_w = e^{-w(z)} \cdot \partial^{-1} \cdot e^{+w(z)} = \left( w'(z) + \partial \right)^{-1}$$

$$\left( \partial = \partial/\partial z ; z = 1/t ; w(z) = \omega_0 + \omega_1 z + \cdots + \omega_q z^q \right)$$

and their equivalents in the various  $\zeta_i$  planes, for all three cases: *precritical* ( $q_i < q$ ), *critical* ( $q_i = q$ ) and *postcritical* ( $q_i > q$ ).

## §4. Acceleration and Quasianalyticity. Cohesive Functions

### Transfinite Denjoy Classes of Quasianalytical Functions. Cohesive Functions

Let  $\mathcal{L}$  be a  $C^\infty$  automorphism of  $]..., +\infty]$  with  $\mathcal{L}(x) \ll x$ . An *iterator*  $\mathcal{L}^*$  of  $\mathcal{L}$  is any  $C^\infty$  automorphism of  $]..., +\infty]$  such that  $\mathcal{L}^* \circ \mathcal{L} = -1 + \mathcal{L}^*$ . For any transfinite ordinal  $\alpha = \omega^r \cdot n_r + \cdots + \omega \cdot n_1 + n_0 < \omega^\omega$  ( $n_i \in \mathbb{N}$ ) we put:

$$(4.1) \quad L_\alpha = (L)^{\circ n_0} \circ (L^*)^{\circ n_1} \circ (L^{**})^{\circ n_2} \circ \cdots \circ (L^{*\cdots*})^{\circ n_r} = \alpha\text{-th iterate of } L = \log.$$

The function  $L_\alpha$  is not uniquely determined (unless  $\alpha < \omega$ ) but the algebra  ${}^\alpha D$  of all  $C^\infty$  functions on  $I = [x_1, x_2]$  with derivatives bounded by:

$$(4.2) \quad |\varphi^{(n)}(x)| < c^n \cdot \left( \frac{L_\alpha(n)}{L'_\alpha(n)} \right)^n \quad (c = c(\varphi) = \text{cst} ; \forall x \in I)$$

depends only on  $\alpha$ . We call it the Denjoy class of order  $\alpha$ . Its elements  $\varphi$  are *quasianalytic*, i.e. they vanish if all their derivatives at a given point vanish. The classes  ${}^\alpha D$  increase with  $\alpha$  and their union for all  $\alpha < \omega^\omega$  is the algebra COHES of *cohesive functions*.

### Weak Accelerates are Cohesive and Cohesive Functions are Weak Accelerates

It can be shown [4, 5] that any cohesive function  $\hat{\phi}_2(\zeta_2)$  on an interval  $[0, \sigma]$  is a *weak accelerate* (i.e. the result of a weak acceleration  $z_1 \rightarrow z_2$  with  $\log z_2 / \log z_1 \rightarrow 1$ ) and, under very mild assumptions on the acceleration, the converse holds: each weak accelerate  $\hat{\phi}_2(\zeta_2)$  is cohesive on some interval  $]0, \sigma]$ , i.e. on each  $[\varepsilon, \sigma]$  with  $\varepsilon > 0$ .

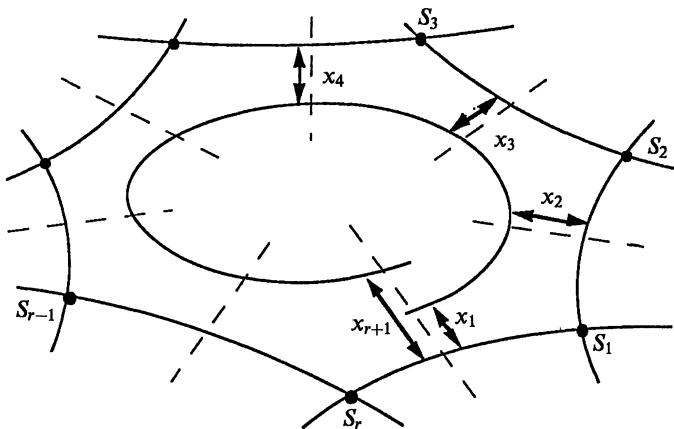
The *cohesiveness* of weak accelerates (just as the *analyticity* of moderate or strong accelerates) is truly providential, for each acceleration  $z_1 \rightarrow z_2$  is actually a two-stepped process: *first*, we calculate  $\hat{\phi}_2(\zeta_2)$  as a *germ* (for small  $\zeta_2$ ) by means of integral (2.6); and *then* we must take its continuation (analytic or quasianalytic) to get  $\hat{\phi}_2$  as a global, multivalued *function* on  $\mathbb{R}^+$ . Of course, when singular points  $\omega_i$  stand in the way of quasianalytic continuation, their “circumvention” (right or left) calls for a special construction [4, 5] since  $\hat{\phi}_2(\zeta_2)$  is not defined outside  $\mathbb{R}^+$ .

The direct statement (cohesive functions are weak accelerates) is also highly meaningful, as it leads to a new and fairly elementary procedure for quasianalytic continuation [4–6].

## §5. The Finiteness of Limit-Cycles. Analysable Functions

I am indebted to J. Martinet, R. Moussu and J.-P. Ramis for drawing my attention to a conjecture by Dulac (long known as Dulac's theorem, but unproven by him) to the effect that the limit-cycles of a vector field on  $\mathbb{R}^2$  with polynomial or real coefficients, cannot possibly accumulate anywhere. Since accumulation could take place only close to a polycycle (or a point) and since polycycles, under repeated blowing-ups, can be brought down to a simple form, the problem may be rephrased as follows. Let  $\mathcal{C}$  be a closed curve on  $\mathbb{R}^2$  consisting of  $r$  analytic arcs  $\mathcal{C}_i$  intersecting at points  $S_i = \mathcal{C}_i \cap \mathcal{C}_{i+1}$ . Let  $X$  be a real analytic vector field defined on a neighbourhood of  $\mathcal{C}$ , with the  $\mathcal{C}_i$  as integral curves and with a non-vanishing linear part at each summit  $S_i$ . Next, draw an analytic curve  $\Gamma_i$  across each  $\mathcal{C}_i$  and endow it with an analytic abscissa  $x_i = 1/z_i$  ( $x_i \sim +0$ ;  $z_i \sim +\infty$ ) positive towards the "interior" of  $\mathcal{C}$ . The integral curve crossing  $\Gamma_i$  at the point with inverse abscissa  $z_i$  crosses  $\Gamma_{i+1}$  at the point  $z_{i+1}$ . The germ  $G_i : z_i \mapsto z_{i+1}$  is the *local map* of summit  $S_i$  and the germ  $F = G_r \circ \cdots \circ G_1$  is the *return map* of  $X$ . Limit-cycles close to  $\mathcal{C}$  clearly correspond to large fixed points of  $F$ . Thus it is all a matter of establishing the trichotomy:

$$(5.1) \quad F(z) \equiv z \quad \text{or} \quad F(z) > z \quad \text{or} \quad F(z) < z \quad (z \gg 1).$$



Due to reduction, the field  $X$  has, at each summit  $S_i$ , either two non-zero eigenvalues of negative ratio  $-\lambda_i \notin \mathbb{Q}$  (type I) or  $-\lambda_i \in \mathbb{Q}$  (type II) or only one non-zero eigenvalue (type III). For all three types, the local map  $G_i$  has a formal counterpart  $\tilde{G}_i$  which is an asymptotic or transasymptotic series of the form:

$$(5.2) \quad \tilde{G}_i = \tilde{K}_i \circ P_{\lambda_i} \circ \tilde{H}_i \quad \text{with } P_{\lambda_i}(z) \equiv z^{\lambda_i} \quad (\text{type I})$$

$$(5.3) \quad \tilde{G}_i = {}^* \tilde{V}_i \circ \tilde{U}_i^* \quad (\text{type II})$$

$$(5.4) \quad \tilde{G}_i = \tilde{K}_i \circ E \circ \tilde{U}_i^* \quad \text{with } E(z) \equiv \exp z \quad (\text{type III}^+)$$

$$(5.5) \quad \tilde{G}_i = {}^* \tilde{V}_i \circ L \circ \tilde{H}_i \quad \text{with } L(z) \equiv \log z \quad (\text{type III}^-)$$

with the  $\tilde{H}_i$  and  $\tilde{K}_i$  denoting ordinary real power series of the form  $az.\{1 + \sum a_n z^{-n}\}$  ( $a > 0$ ,  $a_n \in \mathbb{R}$ ) and with  $\tilde{U}_i^*$  and  ${}^*\tilde{U}_i$  standing for the *formal* iterators (direct and inverse) of the “unitary” maps  $U_i$  which describe the holonomy of  $X$  at  $S_i$ :

$$(5.6) \quad \overline{U} \circ U(z) \equiv z ; \quad \tilde{U}^* \circ {}^*\tilde{U}(z) \equiv z ; \quad \tilde{U}^* \circ \tilde{U}(z) \equiv 2\pi i + \tilde{U}^*(z).$$

The  $\tilde{G}_i$  are usually divergent (except for type I and diophantine  $\lambda_i$ ) but can always be resummed by  $\mathcal{LB}$  with respect to a single critical time  $z_i = h_i(z)$  of the form:

$$(5.7) \quad \begin{cases} h_i(z) = \log z - c \log \log z ; & c \text{ large (type I or II)} \\ h_i(z) \sim \log G_i(z) \text{ (type III$^+$)} & h_i(z) \sim \log z \text{ (type III$^-$)}. \end{cases}$$

As for the return map  $F = G_r \circ \dots \circ G_1$ , its formal counterpart  $\tilde{F} = \tilde{G}_r \circ \dots \circ \tilde{G}_1$  is a transseries with a unique “pulled-down” expansion,

$$(5.8) \quad \tilde{F}(z) = z + \sum a_{\underline{n}} \tilde{A}_{\underline{n}}(z) \quad (0 \leq \underline{n} \leq \underline{n}_0 < \omega^\omega ; a_{\underline{n}} \in \mathbb{R})$$

with finite or transfinite ordinals  $\underline{n}$  and decreasing *transmonomials*  $\tilde{A}_{\underline{n}}$  that are irreducible concatenations of real coefficients and symbols  $+, \times, \log, \exp$ .

Like its factors  $\tilde{G}_i$ , the transseries  $\tilde{F}$  is usually divergent, but always accelero-summable with at most  $r$  critical times  $z_i$  associated with its irregular summits  $S_i$  (actually, the intrinsic notion is that of *critical class*  $\{z_i\}$  regrouping all times equivalent to  $z_i$ ). For each  $z_i$ , only those transseries  $\tilde{\varphi}(z) = \tilde{\varphi}_i(z_i)$  that are *carried* by  $\tilde{F}(z)$  and formally *subexponential* in  $z_i$  possess a Borel transform  $\hat{\varphi}_i(\zeta_i)$  in the  $\zeta_i$  plane. The rest must provisionally retain their status as symbols and bide their “times” to be actualized as true functions! Of course, in order to preserve realness, it is the *median functions*  $\text{med } \hat{\varphi}_i(\zeta_i)$  which are being accelerated or Laplace.

Accelero-summability is proven by induction on  $q$  for  $\tilde{F}_q = \tilde{G}_q \circ \dots \circ \tilde{G}_1$ . Crucial to the argument is the comparability of non-equivalent critical times  $z_i$  and  $z_j$  (one is either faster or slower than the other). Since the function  $F_{ji}$  taking  $z_i$  to  $z_j$  and its transseries  $\tilde{F}_{ji}$  have the same factorization structure as  $F$  and  $\tilde{F}$ , but with a lesser number of factors, one and the same induction takes care of the accelero-summability of  $\tilde{F}$  and  $\tilde{F}_{ji}$ . Despite descriptive and notational hurdles, the proof [4, 5] is amazingly simple\*. It reduces to showing, by induction on  $q$ , that accelero-summing the factors  $\tilde{G}_q$  and  $\tilde{F}_{q-1}$  to  $G_q$  and  $F_{q-1}$  and composing them to get  $F_q = G_q \circ F_{q-1}$ , yields the same result as composing the transseries  $\tilde{G}_q$  and  $\tilde{F}_{q-1}$  to get  $\tilde{F}_q$  and then accelero-summing it to get  $F_q$ . This, in turn, follows from a permutability of type  $\int \sum = \sum \int$ , due to absolute summability, with  $\int$  representing an acceleration or Laplace integral and  $\sum$  standing for the infinite sum which translates, in each model  $\zeta_i$ , the operation of composition  $\circ$ .

\* Another proof, apparently quite different and non-constructive in nature, has been announced by Y.S. Ilyashenko.

We end up with a formal trichotomy:

$$(5.9) \quad \begin{cases} \tilde{F}(z) \equiv z & \text{or} \\ \tilde{F}(z) = z + a_0 \tilde{A}_0(z) + o(\tilde{A}_0(z)) & \text{or} \\ \tilde{F}(z) = z - a_0 \tilde{A}_0(z) + o(\tilde{A}_0(z)) & (a_0 >, \tilde{A}_0 > 0) \end{cases}$$

which, after accelero-summation, translates into the wanted trichotomy (5.1).

### Analysable Functions

The return map  $F$  is only a special instance of *analysable functions*. Those are real-analytic germs on  $]..., +\infty]$  that can be represented by an accelero-summable transseries  $\tilde{F}$  with critical times  $z_i$  which are themselves linked by analysable functions  $F_{ji}$ , etc..., with a *finite critical tree*  $z_{i_1, \dots, i_r}$ . Unlike analytic functions, the class of analysable functions enjoys extreme stability (under all common operations) while retaining the two essential properties of real-analytic functions (*i*) being locally comparable (*ii*) being totally reducible to a formal object, viz. an infinite set of coefficients. Analysable functions are of very frequent occurrence. See [6, 8]. There, we also introduce an even more comprehensive notion of analysable function, which subsumes both complex-analytic and cohesive functions and seems to stretch the Analytic Principle to its farthest possible limit.

### Bibliography

- [1–3] J. Ecalle: Les fonctions résurgentes (vols. 1, 2, 3). Pub. Math. Orsay (1981, 1981, 1985)
- [4] J. Ecalle: Finitude des cycles-limites et accéléro-sommation de l'application de retour. Dans les Actes du Colloque sur les cycles-limites, Luminy 1989. Lecture Notes in Mathematics, vol. 1455. Springer, Berlin Heidelberg New York 1990
- [5] J. Ecalle: La conjecture de Dulac: une preuve constructive. (à paraître à Travaux en Cours, Décembre 1990)
- [6] J. Ecalle: Calcul accélératoire et applications. (à paraître à Travaux en Cours, Décembre 1990)
- [7] J. Ecalle: Calcul compensatoire et linéarisation quasianalytique des champs de vecteurs locaux. (à paraître à Travaux en Cours, 1991)
- [8] J. Ecalle: Fonctions résurgentes, calcul étranger, calcul accélératoire et applications. (cours général – en préparation)
- [9] J. Ecalle, J. Martinet, R. Moussu, J.P. Ramis: Non-accumulation des cycles-limites. C. R. Acad. Sci., série I 304, n° 14 (1298), (I) pp. 375–378, (II) pp. 431–434

# Finiteness Theorems for Limit Cycles

Ju. S. Il'yashenko

Department of Mathematics and Mechanics, Moscow State University  
117234 Moscow, USSR

## 1. Statement of Results

**Theorem 1.** *Polynomial vector field in the real plane has only finitely many limit cycles.*

**Theorem 2.** *Analytic vector field in the closed two dimensional surface has only finitely many limit cycles.*

**Theorem 3.** *A singular point of any analytic vector field has a neighbourhood free of limit cycles.*

**Theorem 4** (Nonaccumulation Theorem). *A polycycle of an analytic vector field in the closed two dimensional surface has a neighbourhood free of limit cycles.*

It is known from the time of Poincaré that the Theorem 4 implies Theorems 1–3. The Theorem 3 is the direct consequence of the Theorem 4: The singular point is a particular case of a polycycle. It is distinguished because the wrong opinion that it was proved long ago is widely spread in mathematical literature and folklore.

In the case when the polycycle is a cycle (contains no singular points) Theorem 4 is an immediate consequence of the analytic dependence of solutions on the initial conditions and a uniqueness theorem for analytic functions, applied to the Poincaré map. In the case of the polycycle (separatrix polygon) an analog of the Poincaré map may be defined (Fig. 1); it is a germ of transformation of the semiinterval onto itself. Its natural form is  $(\mathbb{R}^+, 0) \rightarrow (\mathbb{R}^+, 0)$

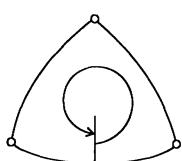


Fig. 1

**Theorem 5 (Identity Theorem).** *If the monodromy map of the polycycle of an analytic vector field has an infinite set of fixpoints, then it is identical.*

Theorem 4 is a trivial consequence of the Theorem 5.

The same results are independently and quite in a different way obtained by J. Ecalle (France).

## 2. Problems Related to the Hilbert's 16th

The classical problem is “... the question on the maximal number and situation of the Poincaré limit cycles for the equation of the form

$$\frac{dy}{dx} = Y/X \quad (1)$$

where  $X, Y$  are polynomials on  $x, y$  of degree  $n$  ...”

The *Hilbert number* for the family of equations (in particular, for a single equation) is the maximum number of limit cycles for the equations of the family. By definition, the Hilbert number is finite or does not exist. Mark some versions of the Hilbert problem.

### Algebraic Versions

1. *The Dulac's problem. Prove the existence of the Hilbert number for the equation (1).*
2. *Prove the existence of the Hilbert number for the whole family (1) (it is noted as  $H(n)$ ).*
3. *Give the (explict) upper estimate on  $H(n)$ .*

### Analytic Versions

4. *Prove the existence of the Hilbert number for the analytic vector field in the sphere  $S^2$ .*
5. *Prove the existence of the Hilbert number for any family of analytic vector fields in  $S^2$  with the finite dimensional compact base.*

### Smooth Version (Hilbert-Arnold Problem)

6. *Prove the existence of the Hilbert number for a “typical” family of smooth vector fields in  $S^2$  with the finite dimensional compact base. Here and below “smooth” is “infinitely smooth”.*

Only the Problems 1 and 4 are solved up to now. There are the following natural implications:

$$\begin{array}{c} 5 \Rightarrow 4 \\ \Downarrow \quad \Downarrow \\ 3 \Rightarrow 2 \Rightarrow 1 \end{array}$$

Remark, that even if the Problem 3 will be solved, it will not imply the solution of the Problems 4–6. The exact meaning of the word “typical” in Problem 6 is the part of the solution. The Hilbert-Arnold problem is the illustration of the heuristical principle: the smooth function behaves like an analytic one, when it is met in the finite dimensional typical family.

### 3. Sources of the Proof

The Dulac's problem concentrates many branches of the theory of differential equations, and the tools, used in its solution, have many other applications. These tools are the following: *desingularization*, *Dulac's asymptotic expansion*, *complex extension in sense of Petrovskii-Landis*, *functional cochains*, *super exact asymptotic series*. The main goal of this talk is to describe these tools. Other applications and further developments are briefly expressed in four Appendices.

### 4. Desingularization

The simplest variant of the desingularization is the polar change of coordinates with the subsequent division by the proper power of the distance to the pasted circle instead of the origin. In more details, let  $v$  be an analytic vector field with an isolated singular point 0. The map  $(r, \varphi) \mapsto (r + 1, \varphi)$ ,  $(r, \varphi)$  being polar coordinates, brings the initial vector field, defined in the neighbourhood of 0 to the vector field, defined in the exterior part of the neighbourhood of the unit circle. This new field may be analytically extended in the interior part of this neighbourhood. After the division by some proper power of  $r - 1$  it becomes a field with a finite number of singular points, located on the unit circle pasted instead of the singular point. These new singular points are in some sense “simpler” than the previous one. It may be repeated several times. The simplest form of singular points obtained by this way is given in the

**Definition.** A singular point of the vector field is called *elementary*, if the linearisation of the field at this point has at least one nonzero eigenvalue.

These eigenvalues are called those of the singular point.

**Bendixson-Seidenberg-Dumortier Theorem.** *After a finite number of steps of the blowing-up process an isolated singular point of an analytic vector field can be split to a finite number of elementary ones.*

Using this theorem one may consider the polycycle in the Identity theorem as to be elementary, that is to say, with only elementary singular points. Other applications and developments of the desingularization method are described in the Appendices I and II.

## 5. The Dulac's Theorem

Dulac (1923) found an asymptotic expansion of the monodromy of the polycycle up to any power of the distance from the polycycle. The more convenient chart is the logarithmic one: if  $x$  is a “natural” chart on the semiinterval  $(R^+, 0)$ , then  $\xi = -\ln x$  is a logarithmic chart.

**Dulac's Theorem.** *Let  $\gamma$  be an elementary polycycle of an analytic vector field in the plane. Then the semitransversal to the polycycle may be so chosen, that the correspondent monodromy map  $\varDelta_\gamma$  will be either flat, or inverse to flat, or admits an asymptotic expansion  $\hat{\varDelta}_\gamma$  having the following form in the logarithmic chart*

$$\hat{\varDelta}_\gamma = \alpha\xi + \beta + \sum P_j(\xi) \exp(-v_j\xi), \quad (2)$$

$\alpha > 0$ ,  $\beta \in R$ ,  $P_j$  are real polynomials,  $0 < v_j \rightarrow \infty$ .

Denote by  $\text{Fix}_\infty$  the set of germs  $(R^+, \infty) \rightarrow (R^+, \infty)$  with the infinite number of fixpoints.

**Dulac's Lemma.** *If  $\varDelta_\gamma \in \text{Fix}_\infty$  then  $\varDelta_\gamma = \text{id}$ .*

The proof, given in (Dulac 1923) uses only the properties of  $\varDelta_\gamma$ , listed in the above theorem. For the maps only with such properties the lemma is wrong; counterexample:  $x \mapsto x + (\exp(-1/x)) \sin 1/x$ . Note, by means of the theorem 5 that such a map cannot appear as a monodromy in the analytic case.

In fact Dulac proves, that the conditions of the lemma imply  $\hat{\varDelta}_\gamma = \text{id}$ . It is trivial, because the terms of the summation in the right hand side of (2) are nonoscillating, and each term tends to zero faster than the previous one.

The Dulac's Theorem is proved by means of smooth, not analytic, theory (Il'yashenko, 1985). The recent development of this theory is described in the Appendix III.

## 6. Complex Extension: The Hyperbolic Case

Let the polycycle  $\gamma$  have only hyperbolic singular points (having no eigenvalues in the imaginary axis). In this case the monodromy map in the logarithmic chart may be extended in the domain, which is “like the right halfplane”  $\mathbb{C}^+ = \{\text{Re } \zeta > 0\}$  and has the form

$$\Omega_C = \Phi_C C^+, \quad \Phi_C : \zeta \mapsto \zeta + C\sqrt{\zeta + 1}$$

On the other hand  $\varDelta_\gamma$  may be decomposed in this domain in the asymptotic series (2). The germs with these two properties are called “almost regular”.

In the hyperbolic case the Identity Theorem may be easily proved (Il'yashenko, 1984). Indeed, if  $\varDelta_\gamma \in \text{Fix}_\infty$ , then  $\hat{\varDelta}_\gamma = \text{id}$  by the Dulac's lemma. Then  $\varDelta_\gamma - \text{id} = o(\exp(-v\xi))$  for any  $v > 0$  on  $(R^+, \infty)$ . At the same time  $\varDelta_\gamma - \text{id}$  is holomorphic and bounded in  $\Omega_C$ . The theorem of Phragmen-Lindelöf type (namely, Watson's Theorem) implies then, that  $\varDelta_\gamma - \text{id} \equiv 0$ . This proves the Identity Theorem in the hyperbolic case. Pass now to the general case.

## 7. Semihyperbolic Singular Points, Functional Cochains and Reduction to Complex Analysis

In the general case the polycycle in the Identity Theorem, assumed to be elementary without loss of generality (see Section 4) contains either hyperbolic, or *semihyperbolic* singular points. These last points have, by definition, one zero and one nonzero eigenvalue. The monodromy of an elementary polycycle may be decomposed in the product

$$\varDelta_\gamma = \varDelta_N \circ \cdots \circ \varDelta_1 \quad (3)$$

of the so called “*correspondence maps*” (Fig. 2) related to the vertices of the polycycle (Fig. 3). Semihyperbolic points bring in the composition (3) the maps of the three following kinds: exponents, logarithms and functional cochains. The appearance of exponents is easily seen in the example  $\dot{x} = x^2$ ,  $\dot{y} = -y$ . The correspondence map of the semiinterval  $x > 0$ ,  $y = e$  onto the semiinterval  $y > 0$ ,  $x = 1$  is  $x \mapsto y(x) = \exp(-1/x)$ . In the logarithmic chart it is  $\xi \mapsto \exp \xi$ .

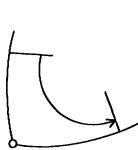


Fig. 2

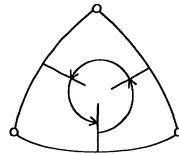


Fig. 3

The functional cochains appear in the description of the correspondence maps for the semihyperbolic singular points of real analytic vector fields. A field of this type has a holomorphic invariant curve, tangent in the semihyperbolic singular point to the eigenvector of the linearization with the nonzero eigenvalue. A monodromy transformation (named the monodromy of the singular point) corresponding to the loop on this curve, surrounding this singular point, is tangent to identity in zero. Such a transformation is formally equivalent to the time one shift along the orbits of the holomorphic vector field with the zero linearization in the singular point. The conjugating formal Taylor series (called also normalizing) is divergent in general. Yet it is asymptotic to some actual function. Namely, there is a covering of a punctured neighbourhood of the fixpoint 0 by sectors of angle  $\alpha \in \left(\frac{\pi}{p}, \frac{2\pi}{p}\right)$

invariant under the rotation by  $\pi/p$  with the vertex at 0. Here  $p + 1$  is the multiplicity of the fixpoint. In each sector a chart is defined, conjugating the initial map with its formal normal form. This chart has an asymptotic expansion at the vertex, given by the normalizing series. The collection of these charts form a normalizing atlas or, in other words, a normalizing cochain. The transition functions of this atlas are called coboundary of the cochain and contain all the information on the geometric properties of the initial germ. They give also the complete invariant of the analytic

classification of the germs  $(\mathbb{C}, 0) \mapsto (\mathbb{C}, 0)$  tangent to identity. The difference between the maps forming the coboundary of the cochain, and identity, is decreasing exponentially and faster, than the nonzero holomorphic function in the sector with the angle, larger than  $\pi/p$ , can decrease.

The correspondence map for the semihyperbolic singular point of a real analytic vector field in the natural map is given by the formula (Il'yashenko, 1986)

$$\Delta = g \circ \exp \circ h_{p,\lambda} \circ H, \quad (4)$$

Here  $H$  is a normalizing cochain for the monodromy of a singular point,  $g : (\mathbb{C}, 0) \mapsto (\mathbb{C}, 0)$  is a holomorphic germ,  $\Delta$  is real on  $(\mathbb{R}^+, 0)$ ,  $h = pz^p/(1 - apz^p \ln z)$

The Identity Theorem is now reduced to the following problem of the complex analysis: prove, that a composition of almost regular germs and of germs, obtained from (4) by transition to the logarithmic map, and also of germs inverse to the previous ones, is either identical, or has no fixpoints near infinity. The ideas for getting a solution of this problem will be described below.

## 8. Phragmen-Lindelöf Property of the Functional Cochains

The class of the functional cochains defined in some domain  $\Omega$  containing  $(\mathbb{R}^+, \infty)$  is called having a *Phragmen-Lindelöf property*, if any cochain of this class, decreasing on  $(\mathbb{R}^+, \infty)$  faster than any exponent:  $\exp(-v\xi)$ ,  $v > 0$ , is identically zero on  $(\mathbb{R}^+, \infty)$ .

All the cochains, appearing in study of the compositions, defined in the end of the Section 7, have the Phragmen-Lindelöf property.

## 9. Super Exact Asymptotic Series

The plan of the further proof is the following: decompose the monodromy map in the asymptotic series with nonoscillating terms; using the Phragmen-Lindelöf Theorem prove, that if this series is equal to identity, then the map itself is identical. The following problem arises in this way. It is not difficult to construct a polycycle of an analytic vector field having the nonidentical monodromy map with the identical Dulac series. Thus the Dulac series does not uniquely determine the monodromy map of the elementary polycycle (though in the hyperbolic case it does). So, we have to construct series, describing both exponential and transexponential decreasing. At first glance it is impossible: any remainder term of the Dulac series may be larger than the transexponential terms. This difficulty can be overcome by the way, shown by the following example. Take two classes of germs of decreasing functions at infinity:  $M_0$  and  $M_1$ , both containing 0 and identity. Let the germs of these classes have the Phragmen-Lindelöf type property: the nonzero germ of the class  $M_0$  ( $M_1$ ) cannot decrease faster than  $\exp(-v\xi)$  (respectively,  $\exp(-v \exp \xi)$ ) for any  $v > 0$ . Let the germs of class  $M_0$  be decomposed in asymptotic Dulac series and the germs of the class  $M_1$  be decomposed in the asymptotic series of the type

$$a_0 + \sum a_j \exp(-v_j \exp \xi) \quad (5)$$

$0 < v_j \rightarrow \infty$ ,  $a_j \in M_0$ ,  $j = 0, 1, \dots$ . The series (5) is called “*super exact asymptotic series*” (the Russian abbreviation is STAR). Its free term contains all the information of the exponential asymptotics of the germ, and the higher terms give the trans-exponential one. The terms of this series do not oscillate and the germs of the class  $M_1$  are uniquely determined by the corresponding STAR’s. Thus for such germs the Identity Theorem holds.

## 10. Additional Decomposition Theorem

In order to apply the previous ideas to the monodromy map, one must decompose it in terms each one of which will be uniquely determined by the corresponding STAR; this STAR is like (5), but much more complicated. The Decomposition Theorem will be stated here for the so called alternant case, when the maps  $\exp$  and  $\ln$  in the composition (3) arise in turn.

**Definition.** Two germs  $f, g$  of functions  $(\mathbb{R}^+, \infty) \mapsto (\mathbb{R}, 0)$  belong to the same class of Archimedean equivalence iff such positive  $a$  and  $b$  exist, that  $|f|^a < |g|$  and  $|g|^b < |f|$  near infinity.

**Additional Decomposition Theorem.** *In the alternant case the monodromy map of the polycycle may be decomposed in the sum*

$$\Delta_\gamma = \alpha\xi + \beta + \varphi + \sum \psi_j \quad (6)$$

with  $|\varphi| < C \exp(-v\xi)$ ,  $|\psi_j| < C_1 \exp(-C_2 \exp v_j \xi)$  for some positive  $v, v_j, C, C_1, C_2$  ( $v_j < v_{j+1}$ ). If  $\alpha = 1, \beta = 0, \varphi \neq 0$ , then  $|\varphi|$  belongs to the same Archimedean class as  $\exp(-\xi)$ ; if  $\alpha = 1, \beta = 0, \varphi \equiv 0$ , then  $|\psi_1|$  belongs to the same Archimedean class as  $\exp(-\exp v_1 \xi)$ .

This theorem implies the Identity Theorem in the alternant case.

**Remark.** In fact, in the decomposition (6)  $\varphi$  has an asymptotic Dulac series, and possesses the Phragmen-Lindelöf property. The germs  $\psi_j \circ v_j^{-1} \circ \ln$  have also the same property; if  $\alpha = 1, \beta = 0, \varphi \equiv 0$ , then the germ  $\psi_1 \circ v_1^{-1}$  may be decomposed in STAR with nonoscillating terms.

## 11. On the Publications and References

The Identity Theorem for the alternant polycycles is proved in full detail in (Il'yashenko, 1990). The complete proof is exposed in a book (Il'yashenko, to appear), which is four times larger than the previous paper forming the first part of the publication. The book contains a complete exposition of a proof with the large introductory chapter, planned for nonspecialists, and is independent of the first part. The first part contains almost all the ideas, used in the book, but is much simpler; it is a kind of digest of the general proof.

The reference list below in relation to Section 7 and Appendices is incomplete for the lack of space. I should mention here that an important contribution in this topic was done by Bogdanov, Bryuno, Denkowska, Dumortier, Ecalle, Elisarov, Malgrange, Martinet, Moussu, Ramis, Roussari, Trifonov Van den Essen, Voronin and others.

## Appendix I. Topological Classification of Nonmonodromic Singular Points in the Plane; the Order of Topologically Sufficient Jet

**Main Alternative.** *A germ of a smooth vector field in the singular point, which is nonflat, is either monodromic (that's to say, admits the Poincaré map, Fig. 4a) or has a characteristic orbit (Fig. 4b). The characteristic orbit is the phase curve, which enters into the singular point after infinite time, in positive or negative sense, tangent to some ray at this point.*

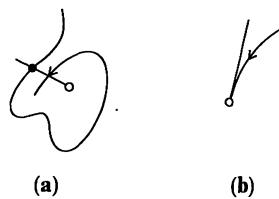


Fig. 4

The Bendixson-Dumortier Theorem guarantees that a finite number of blowing up steps allows to distinguish the two kinds of germs in the Main Alternative, if the singular point has finite multiplicity. By definition, the multiplicity of the singular point of a smooth germ of vector field  $v$  is the dimension of the local ring  $Q = \mathcal{O}/(v)$ , where  $\mathcal{O}$  is the ring of all germs of the smooth functions at the origin, and  $(v)$  is an ideal, generated in this ring by the components of  $v$ . The jet of the vector field is Main-Alternative-sufficient, iff all its representors are simultaneously monodromic or have a characteristic orbit. It is said to be topologically sufficient if all its representors are orbitally topologically equivalent. This means, that the phase curves of any representor may be transformed by the fixpoint preserving homeomorphism into the phase curves of any other representor.

Dumortier proved, that the Main-Alternative-sufficient jet for a smooth germ with finite multiplicity always exists, and for the germ having moreover a characteristic orbit, a topologically sufficient jet also exists.

**Theorem** (Kleban, to appear). *The order of the Main-Alternative- and topologically sufficient jet, defined above, is not larger, than the triple multiplicity of the singular point.*

## Appendix II. Desingularization in the Families of Vector Fields

An important step in understanding of the Hilbert-Arnold's problem is the desingularization theorem for the families of vector (or line) fields. Explain the main statement, omitting the details. A family of complex analytic surfaces is a triple  $M \xrightarrow{\pi} B$  with holomorphic  $n+2$  and  $n$ -manifolds  $M$  and  $B$ , and a holomorphic map  $\pi$  with the constant rank  $n$ . A line field  $\alpha$  on  $M$  is a section of the projective tangent bundle to  $M$  extended in the maximal possible domain  $\Omega_\alpha$  with the following property: the set  $\Sigma = M \setminus \Omega_\alpha$  is analytic and  $\dim \Sigma \leq \dim M - 2$ . By definition, the singular points of the field  $\alpha$  are all the points of  $\Sigma$ . The family of the line fields, corresponding to the family  $M \xrightarrow{\pi} B$  is the line on the total space  $M$ , tangent to the fibers of  $\pi$ . It may be proved, that the restriction of such a family on the fiber, having at least one nonsingular point (such a fiber is called noncritical), may be extended in all the points of the fiber, except some discrete set, called essential singular points. A blowing up of the family of surfaces  $M \xrightarrow{\pi} B$  is the commutative diagram

$$\begin{array}{ccc} \tilde{M} & \xrightarrow{H} & M \\ \downarrow \tilde{\pi} & & \downarrow \pi \\ \tilde{B} & \xrightarrow{\varrho} & B \end{array}$$

where the left column is the family of surfaces with the same dimension of the base,  $H$  and  $\varrho$  are holomorphic, and the restriction of  $H$  on each fiber of  $\tilde{\pi}$  is the finite number of inverse blowing ups. If  $\alpha$  is the family of the vector fields, corresponding to the right column, then the blown up field  $\alpha^*$ , corresponding to the left column, is:  $\alpha = H_* \alpha^*$ .

**Theorem** (Trifonov, 1990). *For any analytic family of line fields, corresponding to the family of surfaces  $M \xrightarrow{\pi} B$  with the compact total space  $M$  with the boundary, having noncritical fibers only, a blowing up exists, giving a new family of line fields, for which all the essential singular points are elementary (that's to say, the extended restriction of the family of fields on each fiber is locally generated by the holomorphic vector field, having only elementary singular points).*

## Appendix III. Smooth Normal Forms in Local Dynamics

The formal equivalence of smooth vector fields is necessary for their smooth equivalence. The inverse implication

$$\text{Formal equivalence} \Rightarrow \text{smooth equivalence} \quad (7)$$

takes place in the following cases:

- 1°. Hyperbolic singular points (Sternberg, Chen).
- 2°. Singular points with one zero eigenvalue (a saddlenode) or an imaginary pair; all other eigenvalues lie outside the imaginary axis (Belitsky, 1986).

For the singular points with the degeneration of the codimension two in the linear part the implication (7) is not proved; in the case of codimension three (for instance, three imaginary pairs) it is wrong (Takens).

Similar results are valid for the local families. In the following cases the families may be brought to the integrable normal forms by a finitely smooth change of coordinates and time (the formulas below give normal forms)

1°.	Perturbation of the nonresonant hyperbolic singular point	$\dot{x} = A(\varepsilon)x$
2°.	Perturbation of the one-resonant hyperbolic singular point (all the resonances are consequences of a single one $(\lambda, r) = 0$ )	$\dot{x} = (\text{diag } x)(\lambda + P(u))$ $u = x^r - a \text{ resonant monomial, } x \in \mathbb{R}^n, P - a \text{ vector polynomial.}$
3°.	Perturbation of a saddle-node of finite multiplicity $p + 1$ with all the resonances being the consequences of $\lambda_1 = 0$	$\dot{x} = x^{p+1}(1 + ax^p)^{-1}$ $\dot{y} = (\text{diag } a(x, \varepsilon))y$ $x \in \mathbb{R}^1, y \in \mathbb{R}^{n-1}$

(Kostov, 1984; Il'yashenko and Yakovenko, 1990)

The finitely smooth classification of the perturbations of the germs with an imaginary pair of eigenvalues has functional moduli (Il'yashenko and Yakovenko, to appear). Parallel theory is developed for diffeomorphisms.

#### Appendix IV. Nonlinear Stokes Phenomena

In the general case local dynamics near a stable point gives a local chart near this point, defined up to a finite number of parameters. The genericity assumption is the Siegel condition, sufficient to the analytic equivalence of the germ of vector field or diffeomorphism to its linear part. The normalizing map is uniquely determined by its linear part. In the resonant case the normalizing chart is replaced by the normalizing atlas, formed by the sector-like domains, covering the punctured neighbourhood of the fixpoint with the point itself on the boundary. The map, conjugating the resonant vector field or the diffeomorphism with its formal normal form is defined and biholomorphic in each sector. The transition functions give the complete invariant of the analytic classification of such germs and contain all the information on the geometric properties of the germ. This program is brought up for germs of maps  $(\mathbb{C}, 0) \rightarrow (\mathbb{C}, 0)$  with the resonant linear part (the multiplicator is the root of unity) and for germs of vector fields in  $(\mathbb{C}^2, 0)$ : the resonant saddles and semihyperbolic singular points (giving rise to the Ecalle-Voronin and Martinet-Ramis moduli). All this material with different extensions and applications is discussed in the forthcoming book (Elisarov et al.).

## References

- Belitsky, G.R.: Smooth equivalence of germs of vector fields with one zero and one nonzero or one pair of pure imaginary eigenvalues. *Funct. Anal. Appl. (Russ.)* **20**, no. 4 (1986) 1–8
- Dulac, H.: Sur les cycles limites. *Bull. Soc. Math. France* **51** (1923) 45–188
- Elisarov, P.M., Il'yashenko, Ju.S., Scherbakov, A.A., Voronin, S.M., Yakovenko, S.Yu.: Nonlinear stokes phenomena. Collection of papers. To appear
- Il'yashenko, Ju.S.: Limit cycles of polynomial vector fields in the real plane with the non-degenerated singular points. *Funct. Anal. Appl. (Russ.)* **18**, no. 3 (1984) 32–42
- Il'yashenko, Ju.S.: The Dulac's memoir “On the limit cycles” and related topics of the local theory of differential equations. *Russ. Math. Surv.* **40**, no. 6 (1985) 41–78
- Il'yashenko, Ju.S.: Separatrix biangles of analytic vector fields in the plane. *Vestnik MGU, Ser. Math.* no. 4 (1986) 25–31
- Il'yashenko, Ju.S.: Finiteness theorems for limit cycles. *Russ. Math. Surv.* **45**, no. 2 (1990) 143–200
- Il'yashenko, Ju.S.: Finiteness theorems for limit cycles. A book, submitted to AMS in spring 1990
- Il'yashenko, Ju.S., Yakovenko, S.Yu.: Normal forms of local families and nonlocal bifurcations. *Russ. Math. Surv.* **46**, no. 1 (1991)
- Il'yashenko, Ju.S., Yakovenko, S.Yu.: Paper V in the book of Elisarov
- Kostov, V.P., Versal deformations of the differential forms of degree  $\alpha$  in the line. *Funct. Anal. Appl. (Russ.)* **18**, no. 4 (1984) 81–82
- Kleban, O.A.: Order of topologically sufficient jet of the vector field in the plane. To appear
- Trifonov, S.I.: Desingularization in the families of analytic vector fields. VINITI, 1990



# Averaging and Passage Through Resonances

Anatoly I. Neishtadt

Space Research Institute, USSR Academy of Sciences, Profsoyuznaya 84/32  
Moscow 117810, USSR

The problems, concerning averaging of perturbations in systems passing through resonances, are dealt with in this report. Such interesting phenomena as capture into resonance, probabilistic scattering of trajectories, destruction of adiabatic invariants, and delay of stability loss, appear to be connected with influence of resonances.

## 1. Slow-Fast Systems, Systems with Rotating Phases, Averaging Method

Slow-fast systems (systems having fast and slow variables) are systems of differential equations of the form

$$\dot{x} = f(x, y, \varepsilon), \quad \dot{y} = \varepsilon g(x, y, \varepsilon), \quad 0 < \varepsilon \ll 1. \quad (1)$$

Variables  $x$  are called fast variables, and variables  $y$  are called slow ones. If  $\varepsilon = 0$ , the system is unperturbed. The equation for  $x$  with  $y = \text{const}$ ,  $\varepsilon = 0$  is called the fast equation (system). The system (1) is investigated using various methods of the perturbation theory corresponding to various properties of the fast equation.

The following case is of great importance for many applications: the fast variables are angle variables (phases) on the  $m$ -torus  $T^m = \mathbb{R}^m / 2\pi\mathbb{Z}^m$ , and a trajectory of the fast system winds round the torus with frequencies  $\omega$  which depend on slow variables. In this case the system (1) is called a system with rotating phases; it has the form

$$\dot{\varphi} = \omega(I) + \varepsilon f(I, \varphi, \varepsilon), \quad \dot{I} = \varepsilon g(I, \varphi, \varepsilon). \quad (2)$$

Below we suppose that right-hand sides of the system are smooth enough (infinitely differentiable for simplicity) functions of all the arguments when  $(\varphi, I, \varepsilon) \in T^m \times D \times [0, \varepsilon_0]$ , where  $D$  is a compact domain in  $\mathbb{R}^n$ . The classical averaging method gives a recipe for approximate description of the slow variable  $I$  evolution in a time interval of order  $1/\varepsilon$ . According to this method one should replace the system (2) with the averaged system

$$\dot{J} = \varepsilon G(J), \quad G(J) = (2\pi)^{-m} \oint_{\mathbb{T}^m} g(J, \varphi, 0) d\varphi. \quad (3)$$

It is supposed that solutions of (3) are good approximations of solutions of (2).

"This principle is neither a theorem, nor an axiom, nor definition, but only a physical suggestion, in other words a vaguely formulated and, to put it strictly, wrong statement. Such statements often appear to be fruitful sources for mathematical theorems" [4].

We denote by  $J(t)$  the solution of the averaged system (3) with initial condition  $I_0 \in D_0 \subset D$  at  $t = 0$  and by  $(I(t), \varphi(t))$  the solution of the slow-fast system (2) (the "precise system") with initial condition  $(I_0, \varphi_0) \in D_0 \times T^m$ . So  $I(0) = J(0) = I_0$ . We suppose that the solution  $J(t)$  is defined and kept at the positive distance from the bound of the domain  $D$  for  $0 \leq t \leq 1/\varepsilon$ . The problem of estimating the difference between  $I(t)$  and  $J(t)$  when  $0 \leq t \leq 1/\varepsilon$  is traditionally called the averaging method justification problem [11].

If the system possesses only one frequency ( $m = 1$ ) and this frequency  $\omega$  does not vanish, the following estimate is valid:  $|I(t) - J(t)| < c\varepsilon$  for  $0 \leq t \leq 1/\varepsilon$ ,  $c = \text{const} > 0$  (P. Fatou; L.I. Mandelshtam and N.D. Papaleksi; this estimate is the first result in the averaging method justification, see [21]). An analogous result holds for multifrequency systems with constant ( $\omega = \text{const}$ ) nonresonant frequencies [11].

For Hamiltonian systems close to integrable ones, the averaging method justification was promoted in frames of the Kolmogorov-Arnold-Moser (KAM) theory [4] and estimates were obtained for an infinite or exponentially large (Nekhoroshev's theorem [4]) time interval. The averaged system in this case is of the form  $\dot{J} = 0$ ; there is no evolution of variables  $I$  and frequencies  $\omega(I)$  in the averaging method approximation.

## 2. Averaging in Multifrequency Systems

In a general case, slow variables  $I$  and, consequently, frequencies  $\omega(I)$  of the system with rotating phases (2) may change in time. In multifrequency ( $m \geq 2$ ) systems at some moments of time a linear dependence of frequencies with integer coefficients, i.e. resonance, occurs. At the resonance a trajectory of the fast system fills a torus of lower dimension and there is no reason to expect that the averaging over the whole torus  $\mathbb{T}^m$  in (3) describes the motion correctly.

The resonance condition  $(k, \omega(I)) = 0$  for each  $k \in \mathbb{Z}^m \setminus \{0\}$  defines a surface, called the resonant surface, in the space of slow variables (here  $(\cdot, \cdot)$  denotes the standard scalar product in  $\mathbb{R}^m$ ). The union of resonant surfaces for all  $k \in \mathbb{Z}^m \setminus \{0\}$  is called the resonant set. The averaging method estimates can be obtained basing on the following idea: if the set of points which are close to resonant surfaces is of small measure, the influence of motion within this set must be small for the majority of initial conditions.

General results were obtained in this direction by D.V. Anosov [1] and T. Kasuga [18]. Let  $E(\varrho, \varepsilon)$  be the set of initial conditions for which the error due to

the averaging method exceeds  $\varrho$  for time  $1/\varepsilon$ :

$$E(\varrho, \varepsilon) = \left\{ (I_0, \varphi_0) \in D_0 \times \mathbb{T}^m : \sup_{0 \leq t \leq 1/\varepsilon} |I(t) - J(t)| > \varrho \quad \text{for} \quad I(0) = J(0) = I_0 \right\} \quad (4)$$

(if  $I(t)$  is not defined on the whole time interval, we formally assume that supremum is equal to 1). According to the theorem by D.V. Anosov [1], if the resonant set is of measure zero, then  $\text{mes } E(\varrho, \varepsilon) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ . Really, the D.V. Anosov theorem has been proved for averaging in general slow-fast systems (1). The following estimate is obtained for system (2).

**Theorem 1** [27]. *Let at least one of the following two conditions be satisfied:*

$$\text{rank}(\partial\omega/\partial I) = m \quad \text{or} \quad \omega \neq 0 \quad \text{and} \quad \text{rank}(\partial(\omega/\|\omega\|)/\partial I) = m - 1. \quad (5)$$

*Then the mean (over initial conditions) error of the averaging method does not exceed a quantity of order  $\sqrt{\varepsilon}$ :*

$$\int_{D_0 \times \mathbb{T}^m} \sup_{0 \leq t \leq 1/\varepsilon} |I(t) - J(t)| dI_0 d\varphi_0 < c_1 \sqrt{\varepsilon}. \quad (6)$$

Here  $c_1$  and furthermore  $c_i$  are positive constants.

**Corollary.** *Under the hypotheses of Theorem 1*

$$\text{mes } E(\varrho, \varepsilon) < c_1 \sqrt{\varepsilon}/\varrho. \quad (7)$$

An equivalent formulation: outside of the set of measure  $\kappa$  the following estimate of the averaging method error is valid:

$$|I(t) - J(t)| < c_1 \sqrt{\varepsilon}/\kappa. \quad (8)$$

The estimate (6) is unimprovable. The estimates (7), (8) are unimprovable within the class of power estimates. It looks reasonable, that these latter estimates can be improved if we confine ourselves to generic perturbations (see below Sect. 3).

If  $m > n + 1$ , the condition (5) cannot be satisfied. In this case the frequency-mapping  $I \mapsto \omega(I)$  defines the submanifold  $M = \omega(D)$  in  $\mathbb{R}^m$ , and to obtain some averaging method estimates one should use Diophantine approximations on this submanifold. The following results have been obtained in this direction.

It is shown in [7] that for arbitrary  $m, n$  the estimates (6)–(8) hold for almost all members of a typical family of frequencies which depends on a sufficiently large number of parameters.

The problem was considered in [16] under the following restriction on the curvature of the manifold  $M = \omega(D)$ . Let  $x \in M$ , vector  $\gamma \in T_x M^\perp$  (the normal space to  $M$  at  $x$ ) and let  $h_\gamma$  be the second fundamental quadratic form of  $M$  with respect to  $\gamma$ . It is assumed, that for every  $x, \gamma$  there exists a two-dimensional subspace in  $T_x M$  where the form  $h_\gamma$  is defined positively or negatively. For  $m < 2 +$

$(n - 1)(n - 2)/2$  this condition is satisfied for generic mappings  $\omega$ , and it is proved that the estimates (6)–(8) are valid.

The estimates for any  $m, n$  and generic mappings  $\omega$  have been obtained.

**Theorem 2** (V.I. Bakhtin [7]). *For systems with  $m$  fast and  $n$  slow variables and with a generic frequency mapping  $\omega$  the mean (over initial conditions) error of the averaging method, does not exceed a quantity of order  $\varepsilon^{1/p+1}$ , provided that  $\binom{n+p}{p} \geq n+m$ .*

Consequently, for such systems

$$\text{mes } E(\varrho, \varepsilon) < c_1 \varepsilon^{1/p+1}/\varrho.$$

The genericity condition is presented in the explicit form. The mappings  $\omega$  not subject to this condition belong to a set of codimension 1 in the functional space.

### 3. Passage Through and Capture into Resonances in Two-Frequency Systems

In two-frequency ( $m = 2$ ,  $\varphi = (\varphi_1, \varphi_2)$ ,  $\omega = (\omega_1, \omega_2)$ ) systems resonant surfaces of different resonances do not cross each other if  $\omega \neq 0$ . So the influence of each resonance can be investigated apart from others.

**Definition [3].** The system (2) meets the condition  $A$  ( $\bar{A}$ , respectively) if the ratio of frequencies  $\omega_1/\omega_2$  changes with a nonzero speed along its trajectories (along the trajectories of the averaged system, respectively):

$$A : (\omega_1 \partial \omega_2 / \partial I - \omega_2 \partial \omega_1 / \partial I)g > c_1^{-1},$$

$$\bar{A} : (\omega_1 \partial \omega_2 / \partial I - \omega_2 \partial \omega_1 / \partial I)G > c_1^{-1},$$

**Theorem 3** (V.I. Arnold [3, 5]). *If the condition  $A$  is satisfied, then*

$$|I(t) - J(t)| < c_2 \sqrt{\varepsilon}, \quad 0 \leq t \leq 1/\varepsilon.$$

**Theorem 4** [24, 26]. *If the system meets both the condition  $\bar{A}$  and some other condition  $B$  (which is actually almost always satisfied), then for all initial points  $(I_0, \varphi_0)$  with the exception of a set  $U_c$  of measure not greater than  $c_2 \sqrt{\varepsilon}$ , the following estimate holds:*

$$|I(t) - J(t)| < c_3 \sqrt{\varepsilon} |\ln \varepsilon|, \quad 0 \leq t \leq 1/\varepsilon.$$

*For any  $\kappa \geq c_2 \sqrt{\varepsilon}$  outside the initial point set of measure not greater than  $\kappa$  the following estimate is valid (compare with (8))*

$$|I(t) - J(t)| < c_4 \sqrt{\varepsilon} |\ln c_5^{-1} \kappa|.$$

The estimates are unimprovable. For the proofs of Theorems 3, 4 see also [20].

For initial conditions belonging to the set  $U_c$  ( $c$  is for “capture”) the capture into resonance takes place. The essence of this phenomenon is that the phase point reaches a resonant surface and begins drifting along this surface, resonance condi-

tions being kept approximately. Therefore, solutions of the precise system and the averaged one diverge by a quantity of order 1 for a time interval  $1/\varepsilon$ . The initial conditions for captured trajectories tend to fill the phase space densely, as  $\varepsilon \rightarrow 0$ . So, it is expedient to consider capture into resonance as a random event and to estimate its probability. At first these phenomena were studied in connection with problems of celestial mechanics [17, 22].

To describe both the capture phenomenon and the behaviour of the captured trajectories, we have to make some additional constructions. It can be shown that for a fixed resonance  $k_1\omega_1 + k_2\omega_2 = 0$  variation of the resonant phase  $\gamma = k_1\varphi_1 + k_2\varphi_2$  close to the resonant surface is described by the following equations:

$$\gamma'' = -\partial V(\gamma, \sigma)/\partial\gamma + L(\sigma) + O(\sqrt{\varepsilon}), \quad \gamma' = O(\sqrt{\varepsilon}). \quad (9)$$

Here  $\sigma \in \mathbb{R}^{n-1}$  are coordinates of the phase point projection onto the resonant surface, and prime denotes the derivative with respect to  $\tau = \sqrt{\varepsilon}t$  [5]. This reduction of the problem was used in [14, 22–24]. Putting in (9)  $\varepsilon = 0$  we obtain the pendulum-like system (shortly, the pendulum) describing the one-dimensional motion in a periodic potential under constant torque  $L$ , where  $L \neq 0$  because of condition  $\bar{A}$ :

$$\gamma'' = -\partial V(\gamma, \sigma)/\partial\gamma + L(\sigma), \quad \sigma = \text{const.} \quad (10)$$

Two possible types of this pendulum's phase portraits are shown in Fig. 1a, b. In Fig. 1a the capture into resonance is impossible. In Fig. 1b there is a domain of oscillation corresponding to phase points captured into resonance. Under the influence of terms  $O(\sqrt{\varepsilon})$  in (9) the phase point can cross the separatrix in Fig. 1b and transit from the rotational domain to the oscillational one. This is a capture into resonance. The backward transition is also possible. This is an escape from resonance.

The oscillational domain of the portrait exists for a finite number of resonances only (because as the order of resonance  $|k| = |k_1| + |k_2|$  grows, the purely periodical term  $\partial V/\partial\gamma$  in (10) tends to 0 and the torque constant  $|L|$  is separated from 0); such resonances are said to be strong. For every oscillational domain of a strong resonance the capture probability can be calculated [31] (this probability can be equal to 0 as well). To describe approximately the motion of captured points within a

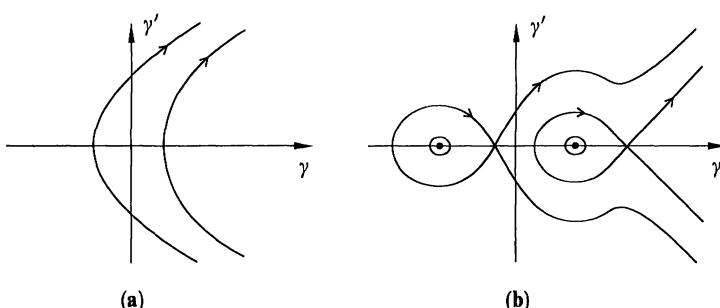


Fig. 1

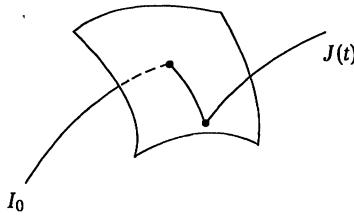


Fig. 2

given oscillational domain we average the velocities of changing of  $\sigma$  and that of the pendulum's energy  $h$  over unperturbed oscillations. So we obtain the set of equations for  $\gamma, h$ , which can be called the internal averaged system for a given oscillational domain of the given resonance. Solving this system one can determine the moment of escape from the resonance if the moment of the capture is known (as the saying is, one can construct the in-out function). To describe approximately trajectories with capture into resonance we glue solutions of the averaged system (3) and the  $\sigma$ -component of solutions of the internal averaged system as it is shown in Fig. 2. It turns out that a finite number of trajectories of the internal averaged system is glued to a trajectory of the averaged system; each of them corresponds to a capture into some oscillational domain of one of strong resonances. It can be shown that under some conditions (satisfied almost always) for the majority of initial conditions in the set  $U_\epsilon$  the behaviour of  $I(t)$  can be described by one of the trajectories glued in this way with an accuracy  $O(\sqrt{\epsilon} \ln \epsilon)$ . The only exception is a subset which measure tends to 0 as  $\epsilon \rightarrow 0$  faster than any given power of  $\epsilon$ .

The condition  $B$  of Theorem 4 is that unstable singular points of the phase portraits are nondegenerate. If the condition  $\bar{A}$  is satisfied (and  $B$  is not), outside a set of initial conditions with measure  $\kappa \geq c_2\sqrt{\epsilon}$  the following unimprovable estimate is valid

$$|I(t) - J(t)| < c_3\sqrt{\epsilon}/\sqrt{\kappa}, \quad 0 \leq t \leq 1/\epsilon$$

(the case  $n = 1$  is considered in [24],  $n \geq 2$ , in [34]).

The estimates in the case where the condition  $\bar{A}$  is not satisfied seem to be unknown. Yet the model problem of estimates has been entirely solved in the case of a single resonance only or, which appears to be the same, for a one-frequency system with the vanishing frequency.

**Theorem 5** (V.I. Bakhtin [8, 9]). *In a generic one-frequency ( $m = 1$ ) system outside a set of initial conditions of measure  $\kappa \geq c_2\sqrt{\epsilon}$  the following estimate holds*

$$|I(t) - J(t)| < c_3\sqrt{\epsilon}/\sqrt{\kappa}, \quad 0 \leq t \leq 1/\epsilon.$$

If  $n \geq 2$  this estimate is unimprovable. The genericity conditions are imposed on the functions  $\omega, G(I)$  near the surface  $\{\omega = 0\}$  and on the function  $\partial\omega/\partial I \cdot g(I, \varphi, 0)$  considered at  $I \in \{\omega = 0\}$ .

In [35] the procedure is proposed, which allows to describe the evolution of slow variables with accuracy better than  $\sqrt{\epsilon}$ , for the case of a single resonance.

Captures into resonance can play the key role in the system's motion for time interval  $\sim 1/\varepsilon^{3/2}$ . Let the averaged system have a periodic solution crossing the resonant surface. Then for time  $1/\varepsilon^{3/2}$  phase points belonging to a set of measure of order 1 will be captured (not to pay attention to changing of the phase volume we deal with systems preserving phase volume). It seems reasonable, that quasi-random captures into resonance and subsequent escapes from resonance can give rise to chaotic dynamics. For examples of such systems see [33, 40].

These phenomena result in destruction of adiabatic invariants in multi-frequency systems (the problem which can be traced back to the paper by P.A.M. Dirac [15]). Let a Hamiltonian system depending on the parameter  $\lambda$  be completely integrable and possess action-angle variables  $I, \varphi$  and frequencies  $\omega = \omega(I, \lambda)$  for each fixed  $\lambda$ . Let  $\lambda$  be a slow periodic function of time:  $\lambda = \lambda(\varepsilon t)$ . Changing of  $I, \varphi$  is described by a system of the form (2), and the averaged system is of the form  $\dot{J} = 0$ , so in the averaging method approximation  $I = \text{const}$ . In one-frequency systems, according to the theorem by V.I. Arnold [2], the action  $I$  remains eternally near its initial value or as the saying goes, appears to be a perpetual adiabatic invariant (provided that some nondegeneracy conditions are satisfied). In [32] an example of a two-frequency system is constructed, where due to captures into resonance  $I$  changes by a quantity of order 1 for a set of initial conditions of measure of order 1 and for time interval  $1/\varepsilon^{3/2}$ . In systems passing through resonance without being captured (for example, under the condition  $A$ ) adiabatic invariants seem to be also destructed, though slowly, for times of order  $1/\varepsilon^2$ .

The approach based on the analysis of joint statistical properties of different passages seems to be fruitful to study multiple passages through resonances. To the author's knowledge, there are no strict results obtained on this way yet.

#### 4. Passage Through a Separatrix

Let us consider a slow-fast system (1) such that the corresponding fast system is a Hamiltonian system with one degree of freedom (on the plane for simplicity). In other words, we consider a system of the form

$$\dot{p} = -\partial E/\partial q + \varepsilon f_1, \quad \dot{q} = \partial E/\partial p + \varepsilon f_2, \quad \dot{z} = \varepsilon f_3, \quad (11)$$

$$E = E(p, q, z), \quad f_i = f_i(p, q, z, \varepsilon), \quad 0 < \varepsilon \ll 1, \quad (p, q) \in \mathbb{R}^2, \quad z \in \mathbb{R}^m$$

with fast variables  $p, q$  and slow variables  $z$ . Let us assume that for all values of  $z$  the phase portrait of the fast system is divided into regions  $G_i$  by separatrices  $l_i$  of a saddle point  $C$  (as, for example, in Fig. 3). A perturbation causes the evolution, and the projection of a phase point onto the  $(p, q)$ -plane may cross the separatrix of the fast system. Far from separatrices the system (11) can be transformed into the form of the system with one rotating phase and a nonzero frequency. To achieve this we have to choose as new variables  $E, z$  and the phase of the unperturbed motion. When phase point approaches the separatrix, the frequency tends to zero. Separatrix may be considered as a special type of resonance.

Passage through a separatrix leads to probability phenomena discovered by I.M. Lifshitz, A.A. Slutskin and V.M. Nabutovskii [19]. Phase points, being initially

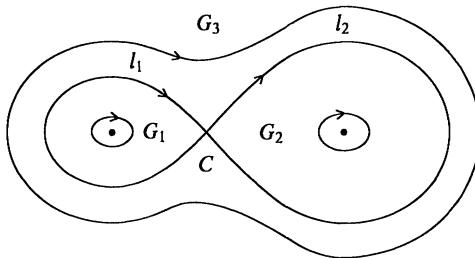


Fig. 3

at the distance of order  $\varepsilon$  from each other, can be captured after separatrix crossing into different regions, and their motions after crossing will be entirely different. Since initial conditions are known always with some finite accuracy, the deterministic approach to the problem fails when  $\varepsilon \rightarrow 0$ . But it is possible to consider the captures into different regions as random events and calculate their probabilities.

Let  $(p_0, q_0) \in G_3$  for  $z = z_0$ . Let  $U_\delta$  be the  $\delta$ -ball around  $M_0 = (p_0, q_0, z_0)$  in  $\mathbb{R}^{n+2}$ . Let  $U_{\delta, \varepsilon}^{(i)}$  be the subset of  $U_\delta$  containing initial points, which will be captured into  $G_i$ ,  $i = 1, 2$ . By definition due to V.I. Arnold [2], the probability of capture into  $G_i$ ,  $i = 1, 2$ , of a point  $M_0$  is

$$Q_i(M_0) = \lim_{\delta \rightarrow 0} \lim_{\varepsilon \rightarrow 0} \frac{\text{mes } U_{\delta, \varepsilon}^{(i)}}{\text{mes } U_\delta}$$

(of course, it must be proved that this limit exists).

Let Hamiltonian  $E = 0$  at the saddle point  $C$  and, therefore, at the separatrices. We will consider the problem under the assumption that the following quantities are different from zero:

$$\Theta_i(z) = - \oint_{l_i} \left( \frac{\partial E}{\partial p} f_1^0 + \frac{\partial E}{\partial q} f_2^0 + \frac{\partial E}{\partial z} f_3^0 \right) dt, \quad i = 1, 2,$$

$$\Theta_3(z) = \Theta_1(z) + \Theta_2(z),$$

the integrals being calculated along the unperturbed separatrices,  $f_j^0 = f_j(p, q, z, 0)$ . These integrals are improper, because the motion along a separatrix takes infinite time. Our normalization of  $E$  ensures the convergence of the integrals. The value  $(-\varepsilon \Theta_j)$  approximates the change of energy  $E$  along the connected part of the perturbed trajectory, which lies near the unperturbed separatrix  $l_j$ . The condition  $\Theta_j \neq 0$  ensures sufficiently fast passage through the separatrix. Let, for certainty,  $\Theta_j > 0$  and  $E > 0$  into  $G_3$ . In this case phase points from  $G_3$  will be captured either into  $G_1$  or into  $G_2$ .

For an approximate description of the evolution we can glue at the separatrices the solutions of the averaged system from different regions. A solution terminating on the separatrix in region  $G_3$  must be glued with solutions beginning at the separatrix in the regions  $G_1$  and  $G_2$  (in Fig. 4 the behaviour of  $E$  for such glued solutions is shown). As is shown in [26], for the majority of initial conditions one

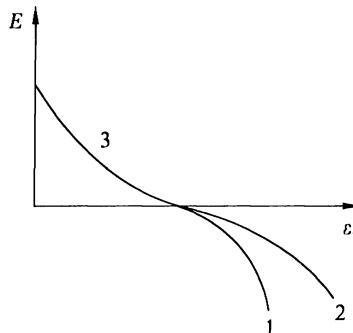


Fig. 4

of these glued solutions describes the evolution of  $E, z$  with accuracy  $O(\varepsilon \ln \varepsilon)$  over time interval of order  $1/\varepsilon$  (an unimprovable estimate). The measure of the set of "bad" initial points, for which such estimate of accuracy is not valid, tends to zero faster than  $\varepsilon^r$  for every given  $r \geq 1$ . The probability of capture into  $G_i$  is calculated by the formula

$$Q_i(M_0) = \frac{\Theta_i(z_*)}{\Theta_1(z_*) + \Theta_2(z_*)}, \quad i = 1, 2 \quad (12)$$

where  $z_*$  is the value of variable  $z$  at the moment of separatrix crossing, calculated by the averaging method. It is easy to understand this formula for probability: the probability of capture into  $G_i$  is equal to ratio of the flux of the phase points in  $G_i$  to the whole flux in  $G_1 \cup G_2$  for  $|z - z_*| < \delta$  in the limit, when  $\varepsilon \rightarrow 0$  and then  $\delta \rightarrow 0$ . Various particular cases of formula (12) were considered in [2, 10, 17, 25].

Analogous phenomena take place and analogous description of the evolution is applicable in the following more general situation. The unperturbed system of differential equations in  $\mathbb{R}^{n+2}$  possesses  $n+1$  first integrals. A joint level of  $n$  integrals is a smooth 2-dimensional manifold. On this manifold the picture of the level lines of the  $(n+1)$ -th integral contains saddle points and separatrices. A perturbation causes the evolution and a phase point crosses these separatrices.

Passages through separatrices destroy adiabatic invariants. Let the system (11) be a Hamiltonian system, depending periodically on slow time  $z = \varepsilon t, f_{1,2} = 0$ . The averaged system possesses an integral (adiabatic invariant) in each region  $G_i$ . This integral is "the action" of the fast system  $I = I(E, z)$  [4]. In the case of separatrix crossing the value of  $I$  along a trajectory changes. Through the period of the slow time this value in general does not return to its initial value even in the averaging method approximation (which is called here the adiabatic approximation). This is an important mechanism of the transport in the phase space of such systems. In addition to this "big" ( $\sim 1$ ) change of  $I$  there exists the "small" ( $\sim \varepsilon |\ln \varepsilon|$ ) change, caused by the difference between exact and averaged solutions. The change of the adiabatic invariant due to the separatrix crossing was calculated through order  $\varepsilon$  first in [38] for the specific case of the simple pendulum in a slowly varying gravity field (the problem considered by P. Ehrenfest) and then in [12, 28] in the general

case of a Hamiltonian system with one degree of freedom, depending slowly on time. The change of the adiabatic invariant strongly depends on initial conditions and must be considered as a random quantity. Apparently, accumulation of these random changes in the case of multiple separatrix crossings leads to diffusion of the adiabatic invariant (this conjecture is confirmed by numerical calculations for many cases). In absence of the separatrix crossing diffusion of the adiabatic invariant does not exist according to Arnold's theorem on perpetual adiabatic invariance [2]. Study of these phenomena requires investigation of the joint statistical properties of different separatrix crossings. Such investigation is now only at the beginning [13].

Similar phenomena take place also in more general Hamiltonian systems with fast and slow motions (the Hamiltonian is  $E = E(p, q, z_1, z_2)$ , where  $(p, q)$  and  $(z_1, \varepsilon^{-1}z_2)$  are pairs of conjugated canonical variables). These phenomena, for example, play an important role in J. Wisdom's theory of the origin of some gaps in the asteroid belt [39, 29].

## 5. Delay of Stability Loss

Let the fast equation of a slow-fast system (1) have an equilibrium position or a limit cycle. Let us suppose that drift of the slow variables leads to a dynamical bifurcation: the equilibrium position (cycle) loses its stability but remains non-degenerate. For an equilibrium a pair of conjugated eigenvalues leaves the left half-plane without passing through point 0, for a cycle multipliers leaves the unit circle without passing through point 1. In analytic systems the stability loss is inevitably delayed [30]: the phase points which moved near the stable equilibrium (cycle) for a time interval of order  $1/\varepsilon$  before the moment of bifurcation, remain near the unstable equilibrium (cycle) for a time of order  $1/\varepsilon$  after the bifurcation. For this time slow variables change by quantity of order 1. Such delay of stability loss is not in general found in nonanalytic systems: receding from the unstable equilibrium takes place near the bifurcation value of the slow variables.

The delay of stability loss for an equilibrium was first discovered by L.S. Pontrjagin and M.A. Shishkova [36] for some model equation system. For this system the asymptotic behaviour of the time of delay was calculated using analytic continuation of solutions in the plane of complex time. The existence of delay itself may be also derived from an earlier statement of Y. Sibuya [37].

The delay of the stability loss for an equilibrium is shown graphically in Fig. 5. It is assumed that the fast motion loses its stability softly, i.e. for  $\tau = \tau_* + 0$  the stable limit cycle is emanating from the equilibrium, and the size of the cycle (the amplitude of selfoscillations) increases like  $\sqrt{\tau - \tau_*}$  (here  $\tau = et$ ). In the precise slow-fast system receding from this equilibrium and transition to the cycle occurs for  $\tau > \tau_* + \text{const}$ , when the cycle is already of size of order 1. The oscillation amplitude increases with a jump (for a time  $t$  of order  $\ln \varepsilon$ ) from the order  $\varepsilon$  to order 1, i.e. the stability loss appears to be produced abruptly.

This phenomenon is in a somewhat unexpected way connected with averaging and passage through resonances. It is convenient to explain this connection using

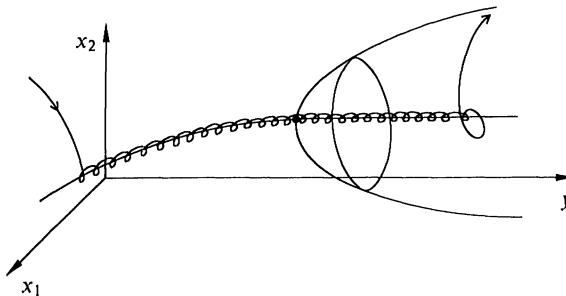


Fig. 5

an example. Let us consider the slow-fast system which is written in a complex form:

$$\dot{z} = (\tau - i)(z - \tau), \quad z = x_1 + ix_2, \quad \tau = \varepsilon t \quad (13)$$

(the equation for the slow variable  $\tau$  has been already integrated, the example in [36] differs from (13) in nonlinear terms). The fast equation has an equilibrium position  $z = \tau$ , which is stable for  $\tau < 0$  and is unstable for  $\tau > 0$ . Let us introduce the function  $\Psi = (\tau - i)^2/2$  and consider, following [36], level lines  $\text{Re } \Psi = \text{const}$  in the plane of the complex slow time  $\tau$  (Fig. 6). If time  $t$  is considered along the path  $\text{Re } \Psi = \text{const}$ , then the equilibrium has purely imaginary eigenvalues. In the polar coordinates  $z = \tau + \varrho e^{i\varphi}$  the system with the rotating phase  $\varphi$  is obtained. The value  $\varrho$  is an integral of the averaged system and an adiabatic invariant of the precise system. At the level lines  $\text{Re } \Psi = 0$ , connecting the points  $-1, i, +1$  (Fig. 6) the resonance occurs: the frequency of  $\varphi$  is equal to zero at the point  $i$ . The passage through this resonance changes the adiabatic invariant by the value of order  $\sqrt{\varepsilon}$ . So it follows that the phase points which were attracted to the equilibrium at the moment of time  $\tau < -1$ , will be ejected at  $\tau \approx +1$ . The phase points, which were attracted at  $\tau = \tau_0 \in [-1, 0]$ , will be ejected at  $\tau \approx |\tau_0|$ .

Analogous phenomena lead to the delay of stability loss of equilibria in the general case. For the motion along some paths in the plane of complex time the system has an adiabatic invariant, which prevents the departure from the equilibrium position, and some critical path with passage through resonance determines the maximal time of delay.

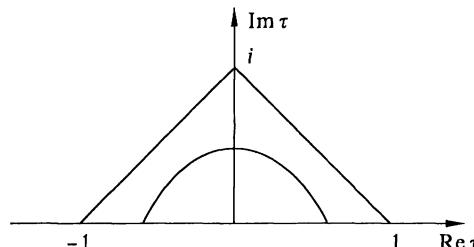


Fig. 6

## References

1. Anosov, D.V.: Averaging in systems of ordinary differential equations with rapidly oscillating solutions. *Izv. Akad. Nauk SSSR, Ser. Mat.* **24** (1960) 721–742 (Russian)
2. Arnold, V.I.: Small denominators and problems of stability of motion in classical and celestial mechanics. *Usp. Mat. Nauk* **18** (1963) 91–192 (Russian). [English transl.: *Russ. Math. Surv.* **18** (1963) 85–192]
3. Arnold, V.I.: Conditions for the applicability and estimate of the error of an averaging method for systems which pass through states of resonance in the course of their evolution. *Dokl. Akad. Nauk SSSR* **161**, 9–12 (1965) (Russian). [English transl.: *Sov. Math. Dokl.* **6** (1965) 331–334]
4. Arnold, V.I.: Mathematical methods of classical mechanics. Nauka, Moscow 1974 (Russian). [English transl.: Springer, New York 1978]
5. Arnold, V.I.: Geometrical methods in the theory of ordinary differential equations. Nauka, Moscow 1978 (Russian). [English transl.: Springer, New York, 1983]
6. Arnold, V.I., Kozlov, V.V., Neishtadt, A.I.: Mathematical aspects of classical and celestial mechanics (Contemporary problems of mathematics. Dynamical systems-3). VINITI, Moscow 1985 (Russian). [English transl.: Encyclopaedia of Mathematical Sciences, vol. 3. Springer, Berlin 1988]
7. Bakhtin, V.I.: Averaging in multifrequency systems. *Funkt. Anal. Prilozh.* **20** (1986) 1–7 (Russian). [English transl.: *Funct. Anal. Appl.* **20** (1986) 83–88]
8. Bakhtin, V.I.: On averaging method in multifrequency systems. Ph.D. Thesis. Moscow State Univ., Moscow, 1986 (Russian)
9. Bakhtin, V.I.: Averaging in generic one-frequency system. *Differ. Uravn.* (to appear) (Russian)
10. Bakai, A.S.: On coefficient of capture of particles into accelerator. *Atomnaya energiya* **21** (1966) 503–504 (Russian)
11. Bogolyubov, N.N., Mitropol'skii Yu. A.: Asymptotic methods in the theory of nonlinear oscillations. 2nd ed. Nauka, Moscow 1958 (Russian) [English transl.: Gordon and Breach, New York 1961]
12. Cary, J.R., Escande, D.F., Tennyson, J.L.: Adiabatic invariant change due to separatrix crossing. *Phys. Rev. A* **34** (1986) 4526–4275
13. Cary J.R., Scodje, R.T.: Phase change between separatrix crossings. *Physica D* **36** (1989) 287–316
14. Chirikov, B.V.: Passage of a nonlinear oscillatory system through resonance. *Dokl. Akad. Nauk SSSR* **125** (1959) 1015–1018 (Russian). [English transl.: *Sov. Phys. Dokl.* **4** (1959) 390–394]
15. Dirac, P.A.M.: The adiabatic invariance of the quantum integral. *Proc. R. Soc. A* **107** (1925) 725–734
16. Dodson, M.M., Rynne, B.P., Vickers, J.A.G.: Averaging in multifrequency systems. *Nonlinearity* **2** (1989) 137–148
17. Goldreich, P., Peale, S.: Spin-orbit coupling in the Solar system. *Astron. J.* **71** (1966) 425–438
18. Kasuga, T.: On the adiabatic theorem for the Hamiltonian system of differential equations in classical mechanics. I, II, III. *Proc. Japan. Acad.* **37** (1961) 366–382
19. Lifshitz, I.M., Slutskin, A.A., Nabutovskii, V.M.: On phenomenon of scattering of charged quasiparticles at singular points in  $p$ -space. *Dokl. Akad. Nauk SSSR* **137** (1961) 553–556 (Russian)
20. Lochak, P., Meunier, C.: Multiphase averaging for classical systems. Springer, New York 1988
21. Mitropol'skii, Yu.A.: Averaging method in nonlinear mechanics. Naukova Dumka, Kiev 1971 (Russian)

22. Moltchanov, A.M.: The resonant structure of the Solar system. *Icarus* **8** (1968) 203–215
23. Morozov, A.D.: A complete qualitative investigation of Duffing's equation. *Differ. Uravn.* **12** (1976) 241–255 (Russian). [English transl.: *Differ. Equations* **12** (1976) 164–174]
24. Neishtadt, A.I.: Passage through a resonance in the two-frequency problem. *Dokl. Akad. Nauk SSSR* **221** (1975) 301–374 (Russian). [English transl.: *Sov. Phys. Dokl.* **20** (1975) 189–191]
25. Neishtadt, A.I.: Passage through a separatrix in a resonance problem with a slowly varying parameter. *Prikl. Mat. Mekh.* **39** (1975) 621–632 (Russian). [English transl.: *J. Appl. Math. Mech.* **39** (1975) 594–605]
26. Neishtadt, A.I.: Some resonance problems in nonlinear systems. Ph.D. Thesis. Moscow State Univ., Moscow 1975 (Russian)
27. Neishtadt, A.I.: Averaging in multifrequency systems. II. *Dokl. Akad. Nauk SSSR* **226** (1976) 1295–1298 (Russian). [English transl.: *Sov. Phys., Dokl.* **21** (1976) 80–82]
28. Neishtadt, A.I.: Change of an adiabatic invariant at a separatrix. *Fiz. Plazmy* **12** (1986) 992–1001 (Russian). [English transl.: *Sov. J. Plasma Phys.* **12** (1986) 568–573]
29. Neishtadt, A.I.: Change of an adiabatic invariant at a separatrix in the systems with two degrees of freedom. *Prikl. Mat. Mekh.* **51** (1987) 750–757 (Russian). [English translation: *PMM USSR* **51** (1987) 586–592]
30. Neishtadt, A.I.: Persistence of stability loss for dynamical bifurcations. I, II. *Differ. Uravn.* **23** (1987) 2060–2067; **24** (1988) 226–233 (Russian). [English transl.: *Differ. Equations* **23** (1987) 1385–1390; **24** (1988) 171–176]
31. Neishtadt, A.I.: Averaging, capture into resonances and chaos in nonlinear systems. In: Campbell, D. (ed.) *Chaos*. AIP, New York 1990, pp. 261–275
32. Neishtadt, A.I.: Probability phenomena in perturbed systems. In: Bazikin A., Zarkhin, Yu. (eds.) *Mathematics and modelling*. Puschino, 1990, pp. 141–155 (Russian)
33. Neishtadt, A.I., Petrovichev, B.A., Chernikov, A.A.: On capture of the particles into regime of unlimited acceleration. *Fiz. Plazmy* **15** (1989) 1026–1029 (Russian). [English translation: *Sov. J. Plasma Phys.* **15** (1989) 593–594]
34. Pronchatov, V.E.: On estimate of the error of the averaging method in two-frequency problem. *Mat. Sb.* **134** (1987) 28–41 (Russian)
35. Sanders, J.: On the passage through resonance. *SIAM J. Math. Anal.* **10** (1979) 1220–1243
36. Shishkova, M.A.: Examination of one system of differential equations with a small parameter in the highest derivatives. *Dokl. Akad. Nauk SSSR* **209** (1973) 576–579 (Russian)
37. Sibuya, Y.: Sur reduction analytique d'un systeme d'équations différentielles ordinaires linéaires content un paramètre. *J. Univ. Tokio, sect. 1*, **10** (1958) 527–540
38. Timofeev, A.V.: On the constancy of an adiabatic invariant when the nature of the motion changes. *Zh. Eksp. Teor. Fiz.* **75** (1978) 1303–1306 (Russian). [English translation: *Sov. Phys. JETP* **48** (1978) 656–659]
39. Widsom, J.: A perturbative treatment of motion near the 3/1 commensurability. *Icarus* **63** (1985) 272–289
40. Zaslavsky, G.M., Neishtadt, A.I., Petrovichev, B.A., Sagdeev, R.Z.: Intensification of diffusion due to interaction wave-particle in a weak magnetic field. *Fiz. Plazmy* **15** (1989) 638–641 (Russian). [English translation: *Sov. J. Plasma Phys.* **15** (1989) 368–370]



# Entropy in Smooth Dynamical Systems

Sheldon E. Newhouse

Department of Mathematics, University of North Carolina  
Chapel Hill, NC 27599-3250, USA

## 1. Introduction

The topological entropy of a continuous dynamical system is now well established as an important invariant of the system. It was first defined by Adler, Konheim, and McAndrew [AKM] in 1965 using open coverings<sup>1</sup>. Nowadays it is convenient to think of the topological entropy as a limit of the number of the distinct orbits of a given length which can be obtained with a fixed small precision. Thus, it measures the *orbit growth* of the system.

More precisely, suppose that  $f : M \rightarrow M$  is a continuous self-mapping of the compact metric space  $M$  with metric  $d$ . Given a positive integer  $n$ , and a small real number  $\delta > 0$ , we say that a set  $E$  is an  $(n, \delta)$ -separated set if, for any  $x \neq y \in E$  there is a  $j \in [0, n)$  such that  $d(f^j x, f^j y) > \delta$ . Let  $r(n, \delta, f)$  be the maximum cardinality of an  $(n, \delta)$ -separated set in  $M$ . Let

$$h(\delta, f) = \limsup_{n \rightarrow \infty} \frac{1}{n} \log r(n, \delta, f)$$

and

$$h(f) = \lim_{\delta \rightarrow 0} h(\delta, f) = \sup_{\delta > 0} h(\delta, f).$$

This is the topological entropy of  $f$ . The most interesting properties of  $h(f)$  arise from its relation to set of invariant probability measures of  $f$ . Let us denote this set by  $\mathcal{M}(f)$ , and recall that it is a compact metrizable set. Given  $\mu \in \mathcal{M}(f)$ , the measure-theoretic or metric entropy,  $h_\mu(f)$ , is defined as follows. For a finite Borel measurable partition  $\alpha$ , let

$$H_\mu(\alpha) = - \sum_{A \in \alpha} \mu(A) \log \mu(A).$$

Given two finite partitions  $\alpha, \beta$ , set  $\alpha \vee \beta = \{A \cap B : A \in \alpha, B \in \beta\}$ . Then, set

$$h_\mu(\alpha, f) = \lim \frac{1}{n} H_\mu\left(\bigvee_{i=0}^{n-1} f^{-i}(\alpha)\right)$$

<sup>1</sup> Adler recently pointed out to us that topological entropy generalizes Shannon's notion of *channel capacity* (see [SW, p.7])

and

$$h_\mu(f) = \sup_{\alpha} h_\mu(\alpha, f) = \lim_{\text{diam } \alpha \rightarrow 0} h_\mu(\alpha, f).$$

The following properties hold:

1.  $h(f^n) = nh(f)$ ,  $h_\mu(f^n) = nh_\mu(f)$  for  $n \in \mathbf{Z}^+$
2.  $h(f^t) = |t| h(f^1)$  if  $\{f^t\}_{t \in \mathbb{R}}$  is a continuous flow.
3.  $h(f) = \sup_{\mu \in \mathcal{M}(f)} h_\mu(f)$
4.  $h(f) = h(g)$  if  $f$  is topologically conjugate to  $g$ ; i.e., there is a homeomorphism  $\phi$  with  $\phi f \phi^{-1} = g$ .
5.  $h(f)$  can be computed using  $(n, \delta)$ -separated subsets of sets of large measure for various  $\mu \in \mathcal{M}(f)$  (see [N1]). That is,

$$h(f) = \sup_{\mu \in \mathcal{M}(f)} \lim_{\sigma \rightarrow 1} \inf_{\mu(A) > \sigma} \lim_{\delta \rightarrow 0} \bar{r}(\delta, A)$$

where

$$\bar{r}(\delta, A) = \limsup_{n \rightarrow \infty} \frac{1}{n} \log \bar{r}(\delta, n, A)$$

and  $\bar{r}(\delta, n, A)$  is the maximal cardinality of an  $(n, \delta)$ -separated subset of  $A$ .

The relationship between topological and metric entropies (statement 3 above) is a combination of the work of Goodman, Goodwyn, and Dinaburg. It is referred to as the Variational Principle for Topological Entropy. An elegant proof has been given by Misiurewicz (see [DGS], pp. 140–146 for a more general result). Statement 4 above states that  $h(f)$  is a topological conjugacy invariant of  $f$ . It is generally not a complete invariant of topological conjugacy. However, for certain important systems it is close to being complete. For instance, Adler and Marcus [AM] have shown that two mixing subshifts of finite type with the same topological entropy are almost topologically conjugate in the sense that each is a boundedly finite to one factor of a third mixing subshift of finite type. Also, Milnor and Thurston [MT] have shown that a piecewise monotone continuous map  $f$  of an interval with positive topological entropy is semi-conjugate in a simple way to a piecewise linear map  $g$  with the same number of turning points as  $f$  and such that the slope of each monotone piece of  $g$  has absolute value equal to  $e^{h(f)}$ .

It was known quite early that some form of regularity affected the finiteness of the topological entropy. Embedding a sequence of larger and larger shift automorphisms topologically in a homeomorphism of the two-sphere shows that homeomorphisms need not have finite topological entropy. However, Kushnirenko proved that every  $C^1$  self-map of a compact manifold has finite topological entropy. A natural question arises: When are there measures  $\mu$  for which  $h_\mu(f) = h(f)$ ? Misiurewicz was the first to construct examples of  $C^k$  diffeomorphisms of compact manifolds with no measure of maximal entropy. His first examples were on manifolds of dimension greater than three, but now it is known that such examples arise in small perturbations of diffeomorphisms having a degenerate homoclinic orbit (an intersection of stable and unstable manifolds of a

hyperbolic saddle point with infinite order contact) even on the two dimensional sphere. His construction worked for every finite  $k$ , but it failed in the  $C^\infty$  case. The question of the existence of measures of maximal entropy was quite important. One consequence of the results described here is that for every  $C^\infty$  self-map  $f$  of a compact smooth manifold, the function  $\mu \rightarrow h_\mu(f)$  is uppersemicontinuous on  $\mathcal{M}(f)$ . Hence,  $f$  does indeed possess measures of maximal entropy.

In this article we shall survey several recent results about topological and metric entropy, particularly as they relate to smooth systems. We view the results described here as part of the natural evolution of certain *topological* aspects of the qualitative theory of dynamical systems following the rich development in the sixties due, in large part, to Anosov, Sinai, and Smale. The recent developments on quadratic mappings due to Carleson and Benedicks may be viewed as part of the evolution of certain *quantitative* aspects of this theory. In the case of uniformly hyperbolic dynamics, these two types of aspects come together beautifully in the theory of *equilibrium states* as described, for example, in [B2]. It is natural to search for a generalization of the Equilibrium State Theory which encompasses all of these results.

## 2. Entropy and Volume

To motivate the results, we first consider the case of mappings of the interval. Let  $M$  be the unit interval and let  $f : M \rightarrow M$  be a continuous map with finitely many turning points. Let  $\log^+(x) = \max(\log x, 0)$ .

Let  $\ell(f) = \text{length of image of } f \text{ with multiplicities}$ :

$$\ell(f) = \int_M |f'(x)| dx$$

**Theorem 1** (Misiurewicz-Szlenk [MS]). *With  $f$  as above, the following results hold.*

1.  $h(f) = \lim_{n \rightarrow \infty} \frac{1}{n} \log^+ \ell(f^n)$ .
2.  $\mu \rightarrow h_\mu(f)$  is uppersemicontinuous on  $\mathcal{M}(f)$ .
3.  $f \rightarrow h(f)$  is uppersemicontinuous for  $f \in C^1$  where one perturbs in the  $C^1$  topology keeping the same number of turning points.
4.  $f \rightarrow h(f)$  is lowersemicontinuous for certain  $f$ .

**Theorem 2** (Misiurewicz [M2]). *The map  $f \rightarrow h(f)$  is lowersemicontinuous for all  $C^0$   $f$ .*

**Corollary 3.** *The map  $f \rightarrow h(f)$  is continuous for  $f$  in the set of  $C^1$  maps with a uniformly bounded number of turning points.*

This last corollary was also proved by Milnor and Thurston [MT].

We consider a direct generalization of the preceding results. A clue as to how to proceed comes from work of Margulis in the sixties on the geodesic flow on compact negatively curved manifolds. Following his work, it became known that

the topological entropy of the time-one map equals the maximum volume growth rate of compact disks in the unstable manifolds.

Let  $D^k$  be the unit  $k$ -disk in  $\mathbf{R}^k$ , and let  $M$  be a  $C^\infty$  manifold. Let  $C^r(M, M)$  be the space of  $C^r$  self maps of  $M$ , and let  $\mathcal{D}^r(M, M)$  be the space of  $C^r$  diffeomorphisms of  $M$  where  $r$  is an integer greater than 1.

A  $C^r$   $k$ -disk in  $M$  is a  $C^r$  map  $\gamma : D^k \rightarrow M$ . Define the  $k$ -volume of  $\gamma$  by

$$|\gamma|_k = \int_{D^k} |\Lambda^k T\gamma| d\lambda$$

where  $d\lambda$  is Lebesgue measure on  $D^k$ , and  $\Lambda^k T\gamma$  is the  $k$ -th exterior power of the derivative  $T\gamma$ .

For  $f \in C^r(M, M)$ , let  $\Lambda$  be a compact  $f$ -invariant set, and let  $U$  be a compact neighborhood of  $\Lambda$ . Given a positive integer  $n$ , set  $W^s(n, U) = \bigcap_{0 \leq j < n} f^{-j}(U)$ . This is just the set of points whose iterates from time 0 through  $n - 1$  remain in  $U$ . For a  $k$ -disk  $\gamma$  in  $U$ , set

$$G_k(\gamma, f, U) = \limsup_{n \rightarrow \infty} \frac{1}{n} \log^+ (|\gamma^{n-1} \circ \gamma| \gamma^{-1}(W^s(n, U))|_k).$$

This is the volume growth of the  $f^{n-1}$ -st iterate of the part of  $\gamma$  which remains in  $U$  from time 0 through time  $n - 1$ .

A collection  $\mathcal{A}$  of  $k$ -disks in  $U$  with  $1 \leq k \leq \dim M$ , is called *ample* for  $\Lambda$  if it contains a subcollection  $\mathcal{A}_1$  for which  $\exists K > 0$  such that

1.  $K^{-1} \leq |D_x \gamma(v)| \leq K$  and  $|D_x^2 \gamma(v, v)| \leq K$   
 $\forall x \in \text{domain } \gamma, |v| = 1$  and  $\gamma \in \mathcal{A}_1$

2. For  $x \in \Lambda$ , and for each  $k$ -dimensional subspace  $H$  of  $T_x M$ , there exists a sequence of  $k$ -disks  $\gamma_1, \gamma_2, \dots \in \mathcal{A}_1$  whose tangent spaces at  $\gamma_i(0)$  approach  $H$  in the Grassmann sense as  $i \rightarrow \infty$ .

If  $M$  is complex analytic with a hermitian metric,  $\Lambda$  and  $U$  are as above, and  $\mathcal{A}$  is a collection of holomorphic disks in  $U$ , then  $\mathcal{A}$  is *holomorphically ample* or *h-ample* if, in condition 2 above,  $H$  is assumed to be a complex subspace of the holomorphic tangent space of  $T_x M$ .

Clearly the family of all disks through points in  $\Lambda$  is ample for  $\Lambda$ . If  $M = \mathbf{R}^N$ , with the usual metric, then the collection  $\mathcal{A}$  of affine disks through points in  $\Lambda$  is ample.

Theorem 4 gives an upper bound for entropy in terms of volume growth rates of smooth disks. An earlier upper bound in terms of the average (relative to Lebesgue measure) of the maximum growth of the norms of the exterior powers of the derivative had been obtained by Przytycki [P] for diffeomorphisms.

**Theorem 4** [N1]. Suppose  $f \in C^r(M, M)$ ,  $\Lambda$  is a compact invariant set,  $U$  is a compact neighborhood of  $\Lambda$ , and  $\mathcal{A}$  is an ample family of  $C^r$  disks for  $\Lambda$ . If  $f$  and  $M$  are complex analytic assume that  $\mathcal{A}$  is  $h$ -ample. Then,

$$h(f | \Lambda) \leq \sup_{\gamma \in \mathcal{A}} G(\gamma, f, U).$$

If  $f \in \mathcal{D}^r(M, M)$ , then

$$h(f | \Lambda) \leq \sup_{\substack{\gamma \in \mathcal{A} \\ \dim \gamma < \dim M}} G(\gamma, f, U).$$

In particular, for  $f \in \mathcal{D}^r(M^2)$ ,  $h(f | \Lambda) \leq \sup_{\text{curves } \gamma} G(\gamma, f, U)$ .

Using Theorem 4, one can give simple proofs of the following results.

**Theorem 5.** 1. Let  $f : \mathbf{R}^N \rightarrow \mathbf{R}^N$  be a polynomial map with coordinate functions of degree  $\leq d$ ; i.e., for  $x \in \mathbf{R}^N$ ,  $f(x) = (f_1(x), f_2(x), \dots, f_N(x))$  with  $f_i(x)$  a polynomial in  $x$  of degree  $\leq d$ . If  $\Lambda$  is a compact invariant set for  $f$ , then

$$h(f | \Lambda) \leq N \log d \quad (\text{Gromov}).$$

2. If  $f : S^2 \rightarrow S^2$  is a rational map of degree  $d$ , then

$$h(f) \leq \log d \quad (\text{Gromov} - \text{Ljubich}).$$

3. If  $f : \mathcal{P}^N(\mathbf{C}) \rightarrow \mathcal{P}^N(\mathbf{C})$  is globally defined and holomorphic, then  $h(f) \leq \log(\text{topological degree})$  (Gromov)

S. Friedland [F] has obtained results on volume growth in quasi-projective varieties which generalize Theorem 5.1.

It is well-known (Misiurewicz-Przytycki [MP]) that if  $M$  is compact and  $f : M \rightarrow M$  is  $C^1$ , then  $h(f) \geq \log(\text{topological degree})$ .

So, Theorems 5.2 and 5.3 above are equalities.

From now on, we assume that  $M$  is a compact  $C^\infty$  manifold.

The above theorems give an upper estimate of  $h(f)$  in terms of volume growth rates of disks. For the lower estimate we have

**Theorem 6** (Yomdin [Y1]). For  $f \in C^\infty(M, M)$ , and any  $C^\infty$  disk  $\gamma$  in  $M$ ,

$$h(f) \geq G(\gamma, f).$$

Yomdin's results can be used to give a proof of a generalization of the Shub entropy conjecture in the  $C^\infty$  case. To recall this conjecture, first define the homology growth of a map  $f$ ,  $HG(f)$ , to be  $\limsup_{n \rightarrow \infty} \frac{1}{n} \log |f_*^n|$  where  $f_* : H_*(M, \mathbf{R}) \rightarrow H_*(M, \mathbf{R})$  is the induced map on the direct sum of the real homology groups of  $f$  (given any norm). The entropy conjecture states that for a  $C^1$  diffeomorphism  $f$  of the compact manifold  $M$ , one has  $h(f) \geq HG(f)$ . Of course,  $HG(f)$  is the same as the maximum logarithm of the absolute values of the eigenvalues of  $f_* : H_*(M, \mathbf{R}) \rightarrow H_*(M, \mathbf{R})$ . Yomdin's results show that this holds for arbitrary  $C^\infty$  maps.

**Corollary 7** (Yomdin) ( $C^\infty$  Entropy Conjecture). *If  $f \in C^\infty(M, M)$ , then  $h(f) \geq HG(f)$ .*

We note that the Entropy Conjecture fails in general for piece-wise linear homeomorphisms although it is true for “typical” piecewise linear maps. There is a large literature on various cases of the entropy conjecture (see [FS]). The general conjecture is still unproved for  $f \in \mathcal{D}^r(M, M)$  with  $1 \leq r < \infty$ .

The next result states that, for positive  $h(f)$ , there always exist disks  $\gamma$  for which  $G(\gamma, f)$  assumes the maximum value. In addition, it can be shown that there are such disks for which  $G(\gamma, f)$  is actually a limit and not just a  $\limsup$ .

**Theorem 8.** *For  $f \in C^\infty(M, M)$ ,  $\mathcal{A}$ , an ample family of  $C^\infty$  disks, we have*

$$h(f) = \sup_{\gamma \in \mathcal{A}} G(\gamma, f) = \max_{\gamma \in \mathcal{A}} G(\gamma, f).$$

In general, the disks  $\gamma$  for which  $G(\gamma, f) = h(f)$  are not easily identifiable. However, for an area decreasing self-embedding of a surface with boundary, the entropy is just the growth rate of the length of the boundary.

**Theorem 9** [N2]. *For  $f \in D^\infty(M^2), \partial M^2 \neq \emptyset$ ,  $f$  area decreasing, we have*

$$h(f) = G(\partial M^2, f).$$

### 3. Continuity Properties of Entropy

The methods used in the proofs of the above results concerning volume growth and entropy have local analogs by which we mean that one considers the growth rates of the cardinalities of  $(n, \delta)$ -separated sets or of the volumes of disks which remain in small neighborhoods of the orbits of given points. These can be combined with general results estimating the defect in uppersemicontinuity of both topological and metric entropy to obtain various continuity properties of entropy for  $C^\infty$  systems.

We begin with a description of the so-called local entropy of a dynamical system  $f : M \rightarrow M$ .

Given  $\Lambda \subset M$ ,  $x \in \Lambda$ ,  $n \in \mathbf{Z}^+$ ,  $\varepsilon > 0$ , set

$$W_x(n, \varepsilon, \Lambda) = \{y \in \Lambda : d(f^j y, f^j x) < \varepsilon \ \forall j \in [0, n]\},$$

$$r(n, \delta, \varepsilon, \Lambda) = \sup_{x \in \Lambda} \max \{ \text{card } E : E \subset W_x(n, \varepsilon, \Lambda), E \text{ is } (n, \delta) - \text{separated} \},$$

$$r(\varepsilon, \Lambda) = \lim_{\delta \rightarrow 0} \limsup_n \frac{1}{n} \log r(n, \delta, \varepsilon, \Lambda).$$

If  $\mu \in \mathcal{M}(f)$ , let

$$h_{\mu \text{loc}}(\varepsilon, f) = \lim_{\sigma \rightarrow 1} \inf_{\mu(\Lambda) > \sigma} r(\varepsilon, \Lambda),$$

and set,

$$h_{\text{loc}}(\varepsilon, f) = \sup_{\mu} h_{\mu \text{loc}}(\varepsilon, f).$$

The quantity  $h_{\text{loc}}(\varepsilon, f)$  is called the  $\varepsilon$ -local entropy of  $f$ , and  $h_{\mu\text{loc}}(\varepsilon, f)$  is called the  $\varepsilon$ -local entropy of  $f$  relative to  $\mu$ .

The next theorem states that  $h_{\text{loc}}(\varepsilon, f)$  gives an upper bound for the difference  $h(f) - h(\varepsilon, f)$  while  $h_{\mu\text{loc}}(\varepsilon, f)$  gives an upper bound for the difference  $h_\mu(f) - h_\mu(\beta, f)$  for any partition  $\beta$  with diameter less than  $\varepsilon$ . Earlier estimates of these differences were given by Bowen in [B].

**Theorem 10** [N2]. *For any continuous self map  $f$  of the compact metric space  $M$ , and  $\varepsilon > 0$ ,*

1.  $h(f) \leq h(\varepsilon, f) + h_{\text{loc}}(\varepsilon, f)$
2. *If  $\mu \in \mathcal{M}(f)$  and  $\beta$  is a finite Borel partition with  $\text{diam } \beta < \varepsilon$ , then*  

$$h_\mu(f) \leq h_\mu(\beta, f) + h_{\mu\text{loc}}(\varepsilon, f)$$

Next we consider local volume growth.

Let  $W_x(n, \varepsilon) = W_x(n, \varepsilon, M)$ .

For a disk  $\gamma$ , set

$$G_{\text{loc}}(\varepsilon, \gamma) = \limsup_{n \rightarrow \infty} \frac{1}{n} \log^+ \sup_{x \in M} |f^{n-1} \circ \gamma| \gamma^{-1}(W_x(n, \varepsilon)) |$$

and

$$G_{\text{loc}}(\varepsilon, f) = \sup_{\gamma} G_{\text{loc}}(\varepsilon, \gamma)$$

**Theorem 11** [N2]. *For  $f \in C^r(M, M)$ ,  $r > 1$ ,*

$$h_{\text{loc}}(\varepsilon, f) \leq G_{\text{loc}}(2\varepsilon, f)$$

**Theorem 12** (Yomdin[Y1]). *For  $f \in C^\infty(M, M)$ ,*

$$\lim_{\varepsilon \rightarrow 0} G_{\text{loc}}(\varepsilon, f) = 0$$

From Theorems 10, 11, 12 with some elementary arguments (see [N2]) we have

**Theorem 13.** 1. *For  $f \in C^\infty(M, M)$ ,  $\mu \rightarrow h_\mu(f)$  is uppersemicontinuous*  
2.  *$f \rightarrow h(f)$  is uppersemicontinuous on  $C^\infty(M, M)$ .*

Yomdin [Y1] has an independent proof of Theorem 13.2.

In general,  $f \rightarrow h(f)$  is not lowersemicontinuous, but Katok has proved that this does hold on surfaces.

**Theorem 14** (Katok).  *$f \rightarrow h(f)$  is lowersemicontinuous on  $\mathcal{D}^2(M^2)$*

**Corollary 15.**  *$f \rightarrow h(f)$  is continuous on  $\mathcal{D}^\infty(M^2)$ .*

**Question 16.** Is the preceding map Holder continuous?

Answer: No

Yomdin [Y2] has examples of curves  $\{f_t\}$  of real analytic maps on  $S^2$ ,  $t \in [-\varepsilon, \varepsilon]$  with  $h(f_t) = 0$  for  $t \in [-\varepsilon, 0]$ , and for  $t \in (0, \varepsilon)$ ,

$$h(f_t) - h(f_0) > C \frac{\log |\log t|}{|\log t|}.$$

In view of Corollary 15, we wish to point out some analogies between the entropy map  $f \rightarrow h(f)$  on  $\mathcal{D}^\infty(M^2)$  and the rotation number map  $f \rightarrow \varrho(f) \in \mathbf{R}/\mathbf{Z}$  on  $Homeo^+(S^1)$ , the set of orientation preserving homeomorphisms of the circle. Both  $h(f)$  and  $\varrho(f)$  are topological invariants which depend continuously on  $f$ . Moreover,  $\varrho(f)$  is rational iff  $f$  has periodic points for  $f \in Homeo^+(S^1)$  while (as proved by Katok)  $h(f)$  is positive iff  $f$  has transverse homoclinic points for  $f \in \mathcal{D}^\infty(M^2)$ .

### Problems

1. (Monotonicity of entropy) For  $f_r(x) = r - x^2$ , the function  $r \rightarrow h(f_r)$  is monotone increasing (Douady and Hubbard). What about  $r \rightarrow h(f_{r,b})$  for fixed  $b$  with  $f_{r,b}(x, y) = (r - x^2 + by, x)$  ?

2. Let  $\mathcal{M}_{\max}(f)$  denote the set of measures of maximal entropy for a mapping  $f$ . As a consequence of Theorem 13, for any  $C^\infty f$ ,  $\mathcal{M}_{\max}(f) \neq \emptyset$ . For  $f \in \mathcal{D}^\infty(M^2)$  with  $h(f) > 0$ , is  $\mathcal{M}_{\max}(f)$  a finite dimensional simplex ?

Related to this problem, Hofbauer ([H]) has developed an interesting theory concerning piecewise monotone mappings of an interval with finitely many monotone continuous pieces. His theory can be described by introducing the notion of a zero-entropy set (0-entropy set). Hofbauer calls these *small* sets.

Let  $f : X \rightarrow X$  be a Borel automorphism of a standard Borel space. A 0-entropy set is an  $f$ -invariant subset  $X_1 \subset X$  such that for any ergodic  $\mu \in \mathcal{M}(f)$ , with  $\mu(X_1) = 1$ , we have  $h_\mu(f) = 0$ . By convention, if  $\mathcal{M}(f) = \emptyset$ , then, every invariant subset of  $X$  is a 0-entropy set. Periodic orbits are simple 0-entropy sets as are stable manifolds of hyperbolic periodic orbits in the smooth setting. There is a natural notion of isomorphism mod 0-entropy: Borel automorphisms  $(f, X), (g, Y)$  are isomorphic mod 0-entropy if there are 0-entropy sets  $X_1 \subset X$ ,  $Y_1 \subset Y$ , and a Borel isomorphism  $\phi : X \setminus X_1 \rightarrow Y \setminus Y_1$ , with  $g\phi = \phi f$ . We say two Borel endomorphisms  $f : X \rightarrow X, g : Y \rightarrow Y$  are isomorphic mod 0-entropy if their natural extensions are isomorphic mod 0-entropy. Finally, we say that the Borel endomorphism  $(f, X)$  is Markov mod 0-entropy if it is isomorphic mod 0-entropy to a finite or countable state topological Markov chain  $(\sigma, \Sigma_A)$ . We will say that a measure preserving endomorphism  $(X, f, \mu)$  is *essentially Markov* if its natural extension  $(\hat{X}, \hat{f}, \hat{\mu})$  is isomorphic to a Markov process. In this case we also say that  $\mu$  is *essentially Markov*.

**Theorem 17** (Hofbauer). *Let  $f : I \rightarrow I$  be a piecewise monotone map of the interval. Then,  $(f, I)$  is Markov mod 0-entropy. Moreover, there are only finitely many ergodic measures of maximal entropy and each is essentially Markov.*

Moving to general  $\dot{C}^\infty$  maps of an interval with positive topological entropy, we can prove that there are at most a countable number of ergodic measures of maximal entropy, and that each is essentially Markov.

In dimension greater than 1, it is not always true that positive entropy implies that  $\mathcal{M}_{\max}(f)$  is a finite dimensional simplex: take the direct product of the identity transformation and any transformation with positive entropy. However, it is possible that, generically, i.e. for elements of a residual set of diffeomorphisms,  $\mathcal{M}_{\max}(f)$  is a finite dimensional simplex.

3. For  $f \in \mathcal{D}^r(M^2)$ ,  $h(f) = 0$ , can  $f$  be  $C^r$  perturbed to be Morse-Smale? This is true if the limit set of  $f$  is finite and hyperbolic [MaP], but not even known if the non-wandering set is finite.

4. For  $f \in \mathcal{D}^\infty(M^2)$ , let  $\chi^+(x)$  denote the positive characteristic exponent of  $x$  (defined for a total probability set of  $x$ ). Let  $\phi(x) = -\chi^+(x)$ . Then,  $\phi$  is bounded and Borel measurable. Define

$$P(\phi) = \sup_{\mu \in \mathcal{M}(f)} h_\mu(f) + \int \phi d\mu.$$

Is there always a  $\mu_0$  such that

$$P(\phi) = h_{\mu_0}(f) + \int \phi d\mu_0.$$

This is true for continuous  $\phi$ . Such  $\mu'_0$ 's are called  $\phi$ -equilibrium states. For a hyperbolic attractor  $A$  and any  $\mu$  supported in the basin of  $A$  and absolutely continuous with respect to Lebesgue measure, it is known that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} f_*^k(\mu)$$

is the unique  $\phi$ -equilibrium state on  $A$  (Ruelle). In general, one would expect  $\phi$ -equilibrium states to be related to weak limits of the measures  $\{\frac{1}{n} \sum_{k=0}^{n-1} f_*^k(\mu)\}$ . In this connection, Pesin and Sinai have shown that the weak limits of the iterates of Lebesgue measure have absolutely continuous densities along unstable leaves for partially hyperbolic attractors [PS].

## References

- [AKM] R. Adler, A. Konheim, M. McAndrew: Topological entropy. Trans. Amer. Math. Soc. **114** (1965) 309–319
- [AM] R. Adler, B. Marcus: Topological entropy and equivalence of dynamical systems. Memoirs Amer. Math. Soc. 219 (1979)
- [B] R. Bowen: Entropy-expansive maps. Trans. Amer. Math. Soc. **164** (1972) 323–331
- [B2] R. Bowen: Equilibrium states and the ergodic theory of Anosov diffeomorphisms. (Lecture Notes in Mathematics, vol. 470.) Springer, Berlin Heidelberg New York 1975

- [DGS] M. Denker, C. Grillenberger, and K. Sigmund: Ergodic Theory on Compact Spaces. (Lecture Notes in Mathematics, vol. 527.) Springer, Berlin Heidelberg New York 1976
- [FS] D. Fried and M. Shub: Entropy, linearity, and chain recurrence. *Publ. Math. IHES* **50** (1979) 203–214
- [F] Friedland: Entropy of polynomial and rational maps. Preprint, Dept. of Math., Univ. of Illinois at Chicago
- [H] F. Hofbauer: On intrinsic ergodicity of piecewise monotonic transformations with positive entropy II. *Israel J. Math.* **38** (1981) no. 1–2, 107–115, and The structure of piecewise monotonic transformations. *Erg. Theory Dyn. Syst.* **1** (1981) no. 2, 159–178
- [M1] M. Misiurewicz: Diffeomorphisms without any measure with maximum entropy. *Bull Acad. Polon. Sci., Ser. Math. Astron. Phys.* **21** (1973) 903–910
- [M2] M. Misiurewicz: Horseshoes for mappings of the interval. *Bull. Acad. Polon. Sci., Ser. Math. Astron. Phys.* **27** (1979) 167–169
- [MaP] I. Malta and M. Pacifico: Breaking cycles on surfaces. *Invent. math.* **74** (1983) no. 1, 43–62.
- [MP] M. Misiurewicz and F. Przytycki: Topological entropy and degree of smooth mappings. *Bull Acad. Polon. Sci., Ser. Math. Astron. Phys.* **25** (1977) no. 6, 573–574
- [MS] M. Misiurewicz and W. Szlenk: Entropy of piecewise monotone mappings. *Asterisque* **50** (1977) 299–310
- [MT] J. Milnor and W. Thurston: Iterated maps of the interval. *Dynamical Systems (Maryland 1986-87)*. (Lecture Notes in Mathematics, vol. 1342) Springer, Berlin Heidelberg New York 1988, pp. 465–563
- [N1] S. Newhouse: Entropy and Volume. *Erg. Theory Dyn. Syst.* **8** (1988) 283–299
- [N2] S. Newhouse: Continuity properties of entropy. *Ann. Math.* **129** (1989) 215–235 and Corrections, *Ann. Math.* **131** (1990) 409–410
- [P] F. Przytycki: An upper estimate for topological entropy of smooth diffeomorphisms. *Invent. math.* **59** (1980) 205–213
- [PS] Pesin-Sinai: Gibbs measures for partially hyperbolic attractors. *Erg. Theory Dyn. Syst.* **2** (1982), no. 3–4, 417–438
- [R] D. Ruelle: A measure associated with Axiom A attractors. *Amer. J. Math.* **98** (1976) 619–654
- [SW] C. Shannon, W. Weaver: The mathematical theory of communication. Univ. of Illinois Press, Urbana, Ill. 1949
- [Y1] Y. Yomdin: Volume growth and entropy. *Israel J. Math.* **57**, no. 3 (1987) 285–301, and  $C^k$ -resolution of semialgebraic mappings – Addendum to ‘Volume growth and entropy’. *Israel J. Math.* **57**, no. 3 (1987) 301–318
- [Y2] Y. Yomdin: Preprint

# Combinatorial Models Illustrating Variation of Dynamics in Families of Rational Maps

Mary Rees

Department of Pure Mathematics, University of Liverpool, P.O. Box 147  
Liverpool L69 3BX, UK

I have been interested in a particular case of the following question:

If  $f_t$  ( $t \in T$ ) is a family of (rational) maps, how do dynamics vary with  $t$ ?

There are a number of advantages in considering this question for families of rational maps. Variation of dynamics in such a family is usually extremely rich, and even the topological nature of the Julia set of a map – the invariant set on which all the interesting dynamics occur – often varies widely. There are usually many hyperbolic maps within such a family. A rational map  $f$  is *hyperbolic* if all critical orbits converge to attractive periodic orbits. In this case, the *Julia set*  $J(f)$  can be defined as the set of all  $z$  such that  $f^n(z)$  does *not* converge to an attractive periodic orbit, and  $f$  is expanding on  $J(f)$ . Such a map  $f$  is relatively easy to analyse, and is dynamically stable in a neighbourhood  $U \supset f^{-1}(U)$  of its Julia set, that is, for all  $g$  sufficiently near  $f$ , there exists a homeomorphism  $\varphi_g : U \rightarrow \overline{\mathbb{C}}$  with  $\varphi_g(J(f)) = J(g)$  and  $\varphi_g \circ f = g \circ \varphi_g$  on  $f^{-1}(U)$ . Thus, most families of rational maps contain both great variation in dynamics and open subsets – stable components – on which dynamics are constant. (It is still not known, however, whether all stable maps are hyperbolic.) Furthermore, a hyperbolic component usually contains a unique *critically finite* map, that is, one for which the critical forward orbits are finite. So dynamics on a hyperbolic component can be examined by examining this map. Finally, perhaps the most important advantage in considering rational families the classical work of Fatou and Julia ([F], [J]) implies that, in general, the dynamics of a rational map is much influenced by the behaviour of its finitely many critical points. This feature, that the dynamics should be rather dependent on the dynamics of finitely many distinguished points, is not unique to families of rational maps.

## A Particular Family

Thanks to the work of Douady and Hubbard [DH1, DH2] and a reinterpretation of some of this by Thurston [T], the best understood family of rational maps is

$$f_a : z \mapsto z^2 + a (a \in \mathbb{C}).$$

The critical point  $\infty$  is fixed throughout this family (as for any polynomial family, of course), and the most significant orbit is that of the other critical point, 0, or of its image, the critical value  $a = f_a(0)$ . Since  $f_a''(\infty) \neq 0$ , there is a holomorphic

map  $\varphi_a$  of some set  $\{z : z > r_a\}$  onto a neighbourhood of  $\infty$  and satisfying  $\varphi_a(z^2) = f_a \circ \varphi_a(z)$ . Moreover,  $a$  is in the image of  $\varphi_a$  for large  $a$ . This very classical result is the key to the very detailed description of the global dynamics of the family  $f_a$ . Douady and Hubbard showed [DH1], [DH2] that the map

$$\Phi : a \mapsto \varphi_a^{-1}(a)$$

is a holomorphic bijection of  $\{a : f_a^n(0) \rightarrow \infty\}$  onto  $\{z : z > 1\}$ , which, significantly, can be rewritten as  $\{z : f_0^n(z) \rightarrow \infty\}$ . Thus, the critical value  $a$  and  $\Phi(a)$  have isomorphic dynamics under the maps  $f_a$ ,  $f_0$  respectively. The map  $\Phi$  extends continuously and with similar isomorphic-dynamics properties to map the Mandelbrot set (see Fig. 1)

$$\{a : f_a^n(0) \not\rightarrow \infty\},$$

not onto the closed unit disc  $\{z : f_0^n(z) \not\rightarrow \infty\}$ , but onto a quotient space of this. This quotient space is interesting, but highly computable (as Thurston's work [T] makes particularly clear). If  $\Phi(a) \leq 1$ ,  $\Phi(a)$  does not always take on under  $f_0$  exactly the dynamics of  $a$  under  $f_a$  (although it often does), partly because  $J(f_a)$  moves with  $a$ . For instance, points on the unit circle (which is the Julia set of  $f_0$ ) come together, and that changes the dynamics of points inside the unit disc. It is not known if  $\Phi$  is a homeomorphism, but recent progress has been made on this outstanding open problem, by Yoccoz.

## Other Families

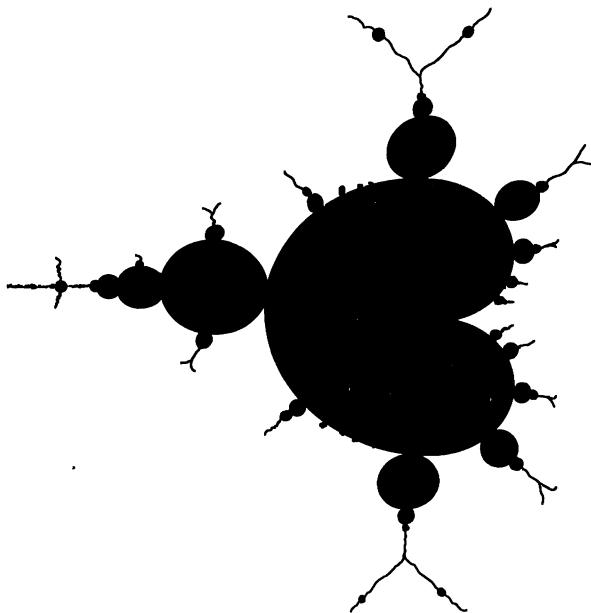
Local dynamics near  $\infty$  are also significant in the study of other polynomial families, as in the work of Branner and Hubbard, for example ([BH1, BH2]). For other families of rational maps, it makes sense to take slices of one complex dimension by keeping the orbits of all but one critical point constant and finite – such slices are conjectured dense (since hyperbolic components are conjectured dense), and one can always examine later the question of how the slices intersect and fit together. A rational map of degree  $d$  has  $2d - 2$  critical points (up to multiplicity), all of multiplicity  $\leq d - 1$ . Thus, a degree two rational map  $f$  has exactly two critical points  $c_1, c_2$ , which vary continuously with  $f$ . I have considered, for  $m \geq 1$ , the variety

$$\{f : f \text{ is degree two, } c_1 \text{ has period } m\}/\text{M\"obius conjugation}$$

If  $m = 1$ , this gives the family  $f_a : z \mapsto z^2 + a$ . For the purposes of discussion, I wish to consider the case  $m = 3$ , which gives

$$g_a : z \mapsto \frac{(z - a)(z - 1)}{z^2} \quad (a \neq 0)$$

in which one critical point 0 has orbit  $0 \mapsto \infty \mapsto 1 \mapsto 0$  of period 3, and the other critical point  $\frac{2a}{a+1}$  has image (or critical value)  $\frac{-(a-1)^2}{4a}$ . There are obviously three values of  $a$  for which this second critical point is fixed, giving, up to M\"obius conjugation, the three polynomials from the family  $\{f_b\}$  with 0 of period 3. Again for the purposes of discussion, we fix one of these polynomials,  $p$ ,



**Fig. 1.** The Mandelbrot set

the anticlockwise rabbit polynomial for which the periodic attractive basins are rotated anticlockwise by  $p$  around the common fixed boundary point. To avoid confusion, (since  $p$  is identified with its Möbius conjugate within the family  $\{g_a\}$ ), we refer to the fixed critical point of  $p$  as  $c_2$ , and the period 3 critical point as  $c_1$ .

## Comparison of a Dynamical Plane and Parameter Space

A crude, inaccurate, but extremely useful idea in analysing such a parameter space is that one should be able to describe the dynamics of any map in the family by movement of one critical value across the dynamical plane of one fixed map in the family, so that the critical value takes in turn the dynamics of all points of the fixed map. We have already indicated how this idea comes into the analysis of the family  $f_a : z \mapsto z^2 + a$ . For the purposes of discussion, we wish, now, to compare the parameter space  $\{g_a : a \neq 0\}$  and the dynamical plane of the anticlockwise rabbit polynomial  $p$ . Within these spaces, the sets

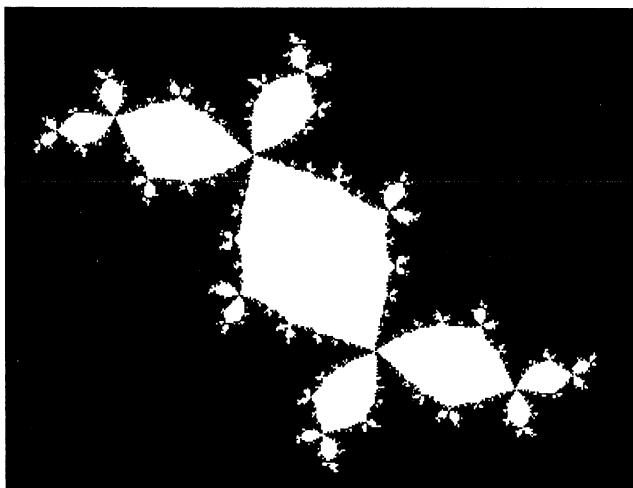
$$\{z : p^n(z) \rightarrow \{c_1, p(c_1), p^2(c_1)\}\}$$

and

$$\{a : g_a^n \left( \frac{2a}{a+1} \right) \rightarrow \{0, 1\infty\}\}$$

should obviously be compared, remembering that 0 and  $\frac{2a}{a+1}$  are the critical points of  $g_a$ . We refer to these sets as the *white sets*, for obvious reasons. See

Figs. 2 and 3. I should like to thank my colleague, F. Rayner, for producing the computer picture on which the parameter space picture is based, and Tan Lei, who wrote the original programme he adapted, and also my colleague R. Morris, for his contribution. Part of the parameter space white set does indeed resemble the dynamical plane white set, and, to a very large extent, this resemblance can be proved. One needs the concept of mating, due to Douady and Hubbard [D], and the concept of capture, which can be found in Wittner's thesis [W]. The complementary black sets are not homeomorphic, though we do not expect that, since there is no such homeomorphism for the family  $\{f_a\}$ . Neither are they clearly dynamically related. So we need, at the very least, to be more precise about the way in which the moving critical value of  $g_a$  can take on the dynamics of points in the dynamical plane of  $p$ .



**Fig. 2.** The dynamical plane of  $p$

## How to Take on the Dynamics of $p$

If we do not try to produce rational maps, it is easy to produce maps taking on the dynamics of various points in the dynamical plane of  $p$ , and this turns out to be worthwhile. We restrict to critically finite branched coverings. If  $f$  is a branched covering, then we define

$$X(f) = \{f^n(c) : c \text{ is critical, } n > 0\}.$$

Then  $f$  is *critically finite* if  $\#(X(f)) < \infty$ . We are going to give some examples based on the critically finite polynomial  $p$  with critical points  $c_1, c_2$ . Let  $\beta : [0, 1] \rightarrow \overline{\mathbb{C}}$  be a path, which for simplicity, we assume is simple, and let  $\beta(0) = c_2$ . Let  $\sigma_\beta$  be a homeomorphism which is the identity outside a (small) disc neighbourhood of  $Image(\beta)$ , and maps  $\beta(0)$  to  $\beta(1)$ . See Fig. 4. Obviously,  $\sigma_\beta \circ p$ , like  $p$ , has critical points  $c_1, c_2$  (provided  $\beta(1) \neq p(c_1)$ ). Provided  $\beta$  avoids the

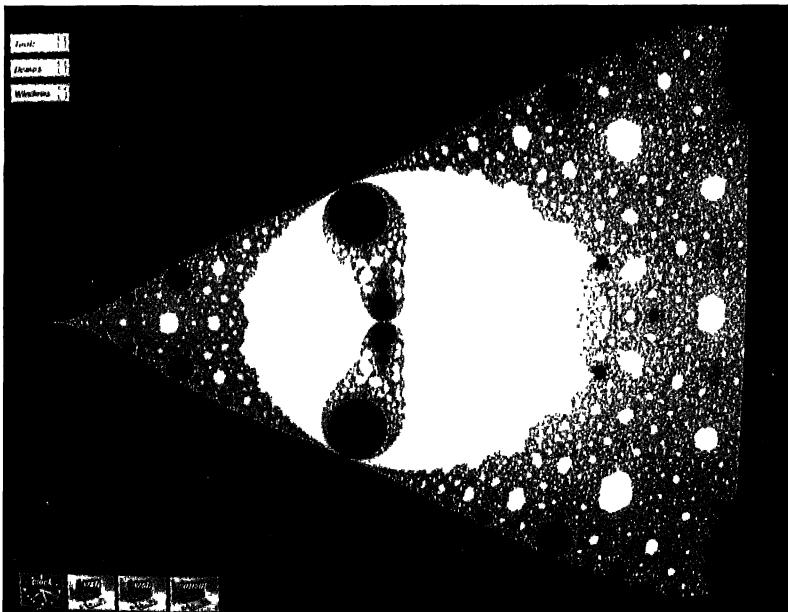


Fig. 3. The parameter space of the  $g_a$

forward orbits of  $\beta(1)$ ,  $c_1$ ,  $\sigma_\beta \circ p$  is a critically finite branched covering, whenever  $\beta(1)$  is a strictly preperiodic point under  $p$ . If  $\beta(1)$  has period  $q > 1$  under  $p$ , then  $\zeta : [0, 1] \rightarrow \overline{\mathbb{C}}$  is uniquely defined by  $p \circ \zeta = \beta$ ,  $\zeta(1) = p^{q-1}(\beta(1))$ , and  $\sigma_\zeta^{-1} \circ \sigma_\beta \circ p$  is a critically finite branched covering, with two periodic critical points.

Obviously, the maps  $\sigma_\beta \circ p$  and  $\sigma_\zeta^{-1} \circ \sigma_\beta \circ p$  are not rational, but they are *equivalent* to rational maps, in a sense yet to be explained. Indeed, we have the following theorem, which, for simplicity, is stated only for the set  $\{g_a : a \neq 0\}$  and the polynomial  $p$ . See [R1].

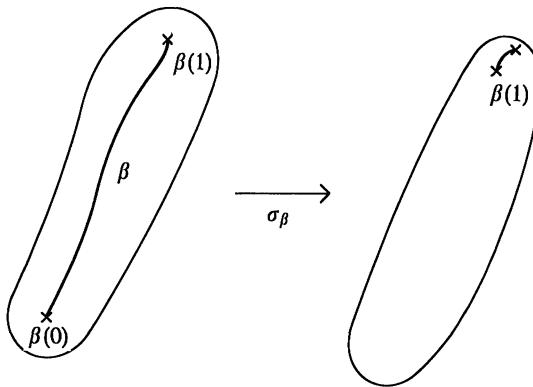
**The Polynomial-and-Path Theorem.** *Any hyperbolic critically finite rational map  $g_a$  is equivalent to  $\sigma_\beta \circ p$  or  $\sigma_\zeta^{-1} \circ \sigma_\beta \circ p$  for some  $\beta$*

In the proof, the path  $\beta$  is obtained from a path  $g_{a_t}$  from  $g_a$  to  $p$ .

## Equivalence of Branched Coverings

This is a homotopy-type equivalence, due to Thurston [T]. The concept occurs in Thurston's theorem [T] about which critically finite branched coverings are equivalent to rational maps. I am avoiding stating this theorem, although it is crucial to the work I am describing.

Let  $f, g$  be critically finite branched coverings of  $\overline{\mathbb{C}}$ . Then  $f$  is *equivalent* to  $g$ , written  $f \simeq g$ , if there exists an orientation-preserving homeomorphism  $\varphi$ , and a path  $g_t$  through critically finite branched coverings, such that  $\varphi \circ f \circ \varphi^{-1} = g_0$ ,  $g = g_1$ ,  $X(g_t) = X(g)$  for all  $t$ .



**Fig. 4.** The homeomorphism

There are many situations in dynamics when some sort of homotopy equivalence between maps  $f$  and  $g$  implies semiconjugacy. Always, some sort of hyperbolicity is required for  $f$ . See, for example, Franks [Fr], where the case of (among others)  $f$  being an Anosov toral automorphism is dealt with. The theorem in the present context is as follows. (See, for example, [R1].)

**The Semiconjugacy Proposition.** *If  $f$  is a critically finite hyperbolic rational map and  $g$  is a critically finite branched covering with  $f \simeq g$ , and  $f$  and  $g$  are conjugate in neighbourhoods of any periodic critical orbits, then there is a continuous map  $\varphi : \overline{\mathbb{C}} \rightarrow \overline{\mathbb{C}}$  such that  $\varphi \circ g = f \circ \varphi$ .*

So a description of a hyperbolic critically finite rational map up to equivalence, as in the Polynomial-and-Path Theorem, is useful. In fact, when the equivalence is to a map of the form  $\sigma_\beta \circ p$  or  $\sigma_\zeta^{-1} \circ \sigma_\beta \circ p$ , the techniques can be extended to describe completely the topological dynamics of the hyperbolic map, and hence of all maps within that hyperbolic component [R1].

## The Information Obtained, and Problems with It

In summary, the Polynomial-and-Path Theorem gives a multivalued map from the set of hyperbolic components within the family  $\{g_a\}$  (the example we are discussing) into a set of paths in the dynamical plane of  $p$ . There are two large problems with this, which we shall refer to as the *Existence Problem* and the *Uniqueness Problem*.

**The Existence Problem.** For what paths in the dynamical plane does there exist a corresponding critically finite rational map? The path space is certainly too large, as it stands (although, of course, no very precise definition has been given).

**The Uniqueness Problem.** How can one decide when two paths give the same rational map?

In order to address these problems, one needs to work with a space obtained by modifying the space of paths. This can be done. See [R1]. (The space is obtained by adapting the ideas of Thurston's combinatorial description of the Mandelbrot set [T], where he also worked out the theory to a considerable extent for other polynomial families. I suspect that similar adaptations may have been made by others, since very substantial work has been done on polynomial families. An adaptation that I am aware of, for a different family, is due to my student D. Ahmadi [A]. Shishikura and Tan Lei [ST] work with a partial combinatorial description in their work on matings of degree 3 polynomials, but following the original Douady-Hubbard approach to the combinatorics.) The modified path space is a *tree with balls*. See Fig. 5 (which is not intended to be accurate). A *tree with balls* is a simply-connected space which is a connected union of a tree and pairwise disjoint balls. In this space, one considers paths from a base point in a ball, and one identifies these paths with their second endpoints.

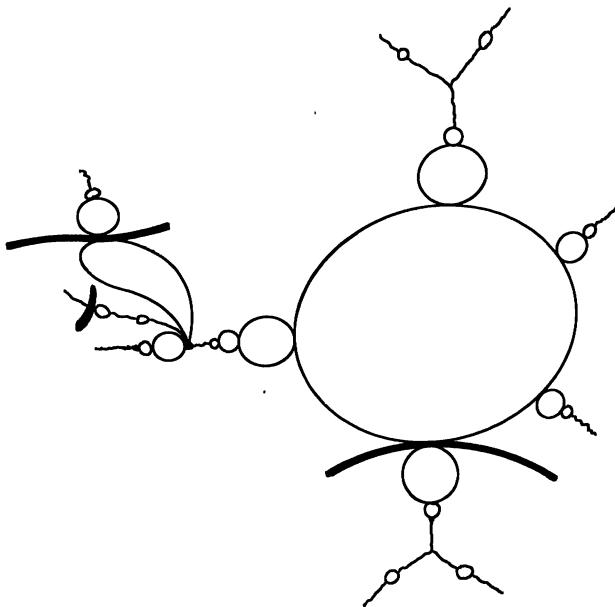


Fig. 5. The tree with balls

All balls in the tree with balls represent critically finite degree two branched coverings (with one critical point of period 3, for the example under discussion). Among them are represented all the hyperbolic critically finite maps with both critical points periodic (in the family  $\{g_a\}$ ), and it makes sense to identify these balls with the corresponding hyperbolic components. Since the tree with balls is a model for the black set in parameter space, boundaries of hyperbolic components from the white set – with one critical point in the backward orbit of the other, periodic, one – are represented simply by points in the tree with balls. It should

be stressed that the tree with balls, and the nature of the balls in it, can be computed entirely from the dynamics of  $p$ , which, in turn, can be computed from the dynamics of  $z \mapsto z^2$ . This is not surprising, since the tree with balls is derived from paths in the dynamical plane of  $p$ . Any hyperbolic component is represented in the tree with balls countably many times, giving countably many paths from the base hyperbolic component. This is rather curious, when there are uncountably many paths between any two hyperbolic components in the true parameter space. Nevertheless, I would hope that the paths in the tree with balls represent *some* of the real paths in parameter space. But little work has been done, as yet, on the relation between the topologies of the two spaces.

## Solution of the Existence Problem

The *Admissible Boundary Theorem* ([R1], [R2]) can be interpreted as follows. Some balls in the tree with balls represent *boundary rational maps*. Their positions, like that of all other balls, can be computed. These boundary balls separate an open set called the *admissible set* from its complement. All balls in the admissible set represent hyperbolic components, and all hyperbolic components are represented there. Possible boundary points for the admissible set (one of which is accurate) are given by the cut-off lines in Fig. 5.

For the family  $\{g_a\}$ , there are two boundary rational maps, which can be associated with the missing points 0 and  $\infty$  in the parameter space. Note that, in fact,  $g_0(z) = \frac{z-1}{z}$ , and after conjugation by  $z \mapsto (-a)^{\frac{1}{2}}z$ ,  $g_a$  converges, as  $a \rightarrow \infty$ , to  $z \mapsto \frac{1}{z}$ . The third and second iterates respectively of these Möbius transformations are the identity. The corresponding *boundary rational maps* are degree two critically finite branched coverings. They are not equivalent to rational maps. Each has a critical point  $c_1$  of period 3. The second critical point  $c_2$  is of period 3 and 2 respectively. In the first case, three circles bounding disjoint discs are cyclically permuted, with one of the discs containing the critical points, one the critical values, and one the images of these. The complement of the discs is invariant under the map. In the second case, there is a simple loop which is mapped to itself, with orientation reversed. It separates the critical values from the other 3 points in the critical forward orbits. See Fig. 6.

## The Uniqueness Problem

One can define an equivalence relation on paths with eventually periodic endpoints in the dynamical plane of  $p$  simply by  $\beta_1 \simeq \beta_2$  if  $\sigma_{\beta_1} \circ p \simeq \sigma_{\beta_2} \circ p$  or  $\sigma_{\zeta_1}^{-1} \circ \sigma_{\beta_1} \circ p \simeq \sigma_{\zeta_2}^{-1} \circ \sigma_{\beta_2} \circ p$  (whichever is appropriate). This equivalence can be transferred to the tree with balls. The equivalence relation on the tree with balls is rather close to being a group orbit equivalence relation.

*Open Question.* Is the equivalence relation sufficiently close to a group orbit equivalence relation for it to be possible to construct a fundamental domain?

Such a fundamental domain would presumably be rather like a Dirichlet fundamental domain.

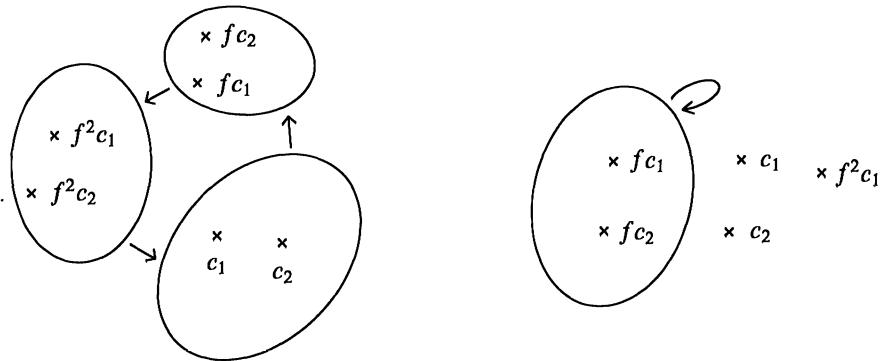


Fig. 6. Boundary rational maps

If  $V = \{g_a : a \neq 0\}$ , then, for  $m \geq 0$ ,

$$V_m = \{g_a : g_a^m \left( \frac{-(a-1)^2}{4a} \right) \notin \{0, 1, \infty\}\}.$$

The techniques of the Polynomial-and-Path Theorem give a natural isomorphism of  $\pi_1(V_m)$  into a certain subgroup  $G(m+1)$  of the group of homeomorphisms of  $\overline{\mathbb{C}}$  preserving  $p^{-m-1}(X(p))$ , modulo the appropriate isotopies, namely, those which are constant on  $p^{-m-1}(X(p))$ . Modulo these, the group  $G(m+1)$  consists of all homeomorphisms  $\varphi$  for which there exists a closed loop  $\alpha$  based at  $c_2$  such that

$$(p, p^{-m-1}(X(p))) \simeq_{\varphi} (\sigma_{\alpha} \circ p, p^{-m-1}(X(p))).$$

This notation means that there is a path  $g_t$  through critically finite branched coverings such that  $g_0 = \varphi \circ p \circ \varphi^{-1}$ ,  $g_1 = \sigma_{\alpha} \circ p$ ,  $g_t^{-m-1}(X(g_t))$  is independent of  $t$ . (The relation  $\simeq_{\varphi}$  is only an equivalence relation when  $\varphi$  is allowed to vary.) The isomorphism is not surjective for  $m \geq 1$ . This seems to be because  $V$  sits in a natural larger space, which includes nonrational maps.

## References

- [A] Ahmadi, D.: Thesis in preparation.
- [BH1] Branner, B., Hubbard, J.H.: The iteration of cubic polynomials. Part I: The global topology of parameter space. *Acta Math.* **160** (1988) 143–206
- [BH2] Branner, B., Hubbard, J.H.: The iteration of cubic polynomials. Part II: Patterns and parapatterns. *Danmarks Tekniske Højskole*, 1989 (preprint)
- [D] Douady, A.: Systèmes dynamiques holomorphes. Séminaire Bourbaki 1983, *Astérisque* **105/6** (1983) 39–63
- [DH1] Douady, A., Hubbard, J.H.: Etudes dynamiques des polynômes complexes, avec la collaboration de P. Lavaurs, Tan Lei, P. Sentenac. Parts I and II. *Publications Mathématiques d'Orsay*, 1985
- [DH2] Douady, A., Hubbard, J.H.: Itérations des polynômes quadratiques complexes. *C.R. Acad. Sci. Paris, Série I* **294** (1982) 123–126

- [F] Fatou, P. : Mémoire sur les équations fonctionnelles. *Bull. Soc. Math. France* **47** (1919) 161–271, **48** (1920) 33–96 and 208–314
- [Fr] Franks, J.: Anosov diffeomorphisms. *Proc. Symp. Pure Math.* **14** (1968) 61–93
- [J] Julia, G.: Itérations des applications fonctionnelles. *J. Math. Pures Appl.* **8** (1918) 47–245
- [R1] Rees, M.: A partial description of parameter space of rational maps of degree two: Part 1. To appear in *Acta Math.*
- [R2] Rees, M.: A partial description of parameter space of rational maps of degree two: Part 2. University of Liverpool (preprint)
- [ST] Shishikura, M., Tan Lei: A family of rational maps and matings of cubic polynomials. Max Planck Institute, Bonn, 1988 (preprint)
- [T] Thurston, W.P.: On the combinatorics of iterated rational maps. Princeton University and I.A.S., 1985 (preprint)
- [W] Wittner, B.: On the bifurcation loci of rational maps of degree two. Thesis, Cornell University, 1988

# Hyperelliptic Function Parametrization for the Chiral Potts Model

R. J. Baxter

Mathematics and Theoretical Physics, I.A.S., The Australian National University  
G.P.O. Box 4, Canberra, A.C.T. 2601, Australia

**Abstract.** The chiral Potts model in statistical mechanics is characterized by sets of variables  $a, b, c, d$ , lying on a high-degree algebraic curve. (Each such set being a “rapidity” vector.) Here we show how they can be parametrized by the hyperelliptic theta functions used by Sonya Kowalevski.

## Introduction

Considerable progress has been made in the last three years in understanding the integrable chiral Potts model [1]–[27]. This statistical mechanical model was first formulated [4, 10] in terms of homogeneous sets of variables  $(a, b, c, d)$  satisfying the relations (only two of which are independent)

$$\begin{aligned} a^N + k'b^N &= kd^N, \quad k'a^N + b^N = kc^N, \\ ka^N + k'c^N &= d^N, \quad kb^N + k'd^N = c^N. \end{aligned} \tag{1}$$

Here  $N$  is an integer greater than one and  $k, k'$  are real constants satisfying  $k^2 + k'^2 = 1$ : it is convenient to introduce a parameter  $\theta$  such that

$$k = \sin \theta, \quad k' = \cos \theta, \quad 0 < \theta < \pi/2. \tag{2}$$

For  $N = 2$ , the relations (1) are of genus 1 and can be uniformized using Jacobi’s elliptic functions. For  $N > 2$  the relations are of higher genus and one needs hyperelliptic functions. Here we show how this can be done.

It has been found convenient to introduce [5, 20] other variables  $v, u, x, y, \mu$ , related to  $a : b : c : d$  and satisfying

$$k \sin u = \sin v, \tag{3a}$$

$$x = e^{i(u-v)/N}, \quad y = e^{i(\pi+u+v)/N}, \tag{3b}$$

$$\mu^N = k'/(1 - kx^N) = (1 - ky^N)/k', \tag{3c}$$

$$a : b : c : d = x\mu : y : 1 : \mu. \tag{3d}$$

Once  $N, \theta$  and any one of the variables  $v, u, x, y, \mu$  is given, the rest are determined (to within a discrete set of choices). Here we regard  $v$  as an independent complex variable;  $u, x, y, \mu$ , and the ratios  $a : b : c : d$ , as functions thereof.

We have shown [27] that low-temperature corner transfer matrix calculations introduce integrals of the form

$$\int^v \frac{e^{i(N-2\alpha)v'/N} dv'}{\Delta(v')}, \quad (4)$$

where  $\alpha = 1, \dots, N-1$  and

$$\Delta(v) = -ik \cos u = \sqrt{\{\sin^2 v - \sin^2 \theta\}}. \quad (5)$$

Although these calculations were only for  $k'$  small, these are *precisely* the hyperelliptic integrals that occur if we specialize the classic work of Kowalevski [28] to the function  $\Delta(v)$ , considering it as the square root of a Laurent polynomial in  $e^{2iv/N}$  (or equivalently in  $\cot[(v + \pi - \theta)/N]$ ).

Here we make this specialization and find, rather remarkably, that it provides a uniformizing parametrization, not just of  $\Delta$  and  $e^{2iv/N}$ , but of all the variables  $v, u, \dots, d$ . Indeed, our main result (48) is that  $a, b, c, d$  can be normalized so that (to within elementary factors) each is an hyperelliptic theta function.

We hope this parametrization will be useful in pursuing the corner transfer matrix calculations.

## Kowalevski's Notation

In §6 (pp. 217–221) of [28], take  $\varrho = N - 1$ ,  $A_0 < 0$ , and define  $a_0, \dots, a_{2\varrho}$  by

$$a_{2j-1} = -\cot(\pi j/N), \quad a_{2j} = -\cot[(\pi j + \pi - 2\theta)/N]. \quad (6)$$

Then  $a_0 < a_1 < \dots < a_{2\varrho}$ . Changing the variable  $x$  in [28] to  $v$ , where  $x = -\cot[(v + \pi - \theta)/N]$ , the expression  $\sqrt{R(x)}$  therein becomes

$$\sqrt{R(x)} = (-1)^{N-1} c_0 \Delta(v) / \{\sin[(v + \pi - \theta)/N]\}^N, \quad (7)$$

where  $c_0 = [-A_0/(N \sin 2\theta)]^{1/2}$  is a positive real constant. The values of  $v$  corresponding to  $x = a_0, \dots, a_{2\varrho}$  are  $v = b_0, \dots, b_{2\varrho}$ , where

$$b_{2j-1} = \pi j - \pi + \theta, \quad b_{2j} = \pi j - \theta.$$

The sign conventions in [28] are such that  $i^{2N-j} \sqrt{R(x)}$  is real and positive for  $a_{j-1} < x < a_j$ . It follows that  $\Delta(v)$  is negative imaginary for  $-\theta < v < \theta$ , positive real for  $\theta < v < \pi - \theta$ , and satisfies the anti-periodicity relation

$$\Delta(v + \pi) = -\Delta(v). \quad (8)$$

The continuation of  $\Delta(v)$  is therefore analytic in the lower half of the complex  $v$ -plane. When  $\text{Im}(v) \rightarrow -\infty$ , then  $\Delta(v) \rightarrow -i e^{iv}/2$ .

The functions  $F_1(x), \dots, F_{N-1}(x)$  of [28] are linearly independent polynomials of degree  $N - 2$ . Allowing complex coefficients, we can choose them so that, for  $\alpha = 1, \dots, N - 1$

$$\frac{F_\alpha(x)}{\sqrt{R(x)}} dx = \frac{e^{i(N-2\alpha)v/N}}{\Delta(v)} dv \quad (9)$$

which is proportional to the integrand that occurs in (4).

## Hyperelliptic Integrals

With these notations, the definite integrals  $K_{\alpha\beta}$ ,  $\bar{K}_{\alpha\beta}$  of [28] are (for  $\alpha, \beta = 1, \dots, N - 1$ )

$$\begin{aligned} K_{\alpha\beta} &= \int_{\pi\beta-\pi+\theta}^{\pi\beta-\theta} \frac{e^{i(N-2\alpha)v/N}}{\Delta(v)} dv, \\ \bar{K}_{\alpha\beta} &= -i \int_{\pi\beta-\pi-\theta}^{\pi\beta-\pi+\theta} \frac{e^{i(N-2\alpha)v/N}}{\Delta(v)} dv. \end{aligned} \quad (10)$$

In [27] we defined the function

$$I(\theta, \alpha) = \int_{-\theta}^{\theta} \frac{e^{i(1-2\alpha)y}}{\sqrt{\{\sin^2 \theta - \sin^2 y\}}} dy \quad (11)$$

and showed that  $I(\theta, \alpha) = \pi F(\alpha, 1 - \alpha; 1; \sin^2 \theta)$ , where  $F(\alpha, \beta; \gamma; z)$  is the usual hypergeometric function. Define

$$L_\alpha = I\left(\frac{\pi}{2} - \theta, \alpha/N\right), \quad L'_\alpha = I(\theta, \alpha/N). \quad (12)$$

Then  $L_\alpha$  and  $L'_\alpha$  are positive real and

$$L_{N-\alpha} = L_\alpha, \quad L'_{N-\alpha} = L'_\alpha \quad (13)$$

$$K_{\alpha\beta} = i\omega^{-\alpha\beta} e^{\pi i \alpha / N} L_\alpha, \quad \bar{K}_{\alpha\beta} = \omega^{\alpha-\alpha\beta} L'_\alpha, \quad (14)$$

where  $\omega = e^{2\pi i / N}$ .

The  $G_{\alpha\beta}$  defined in [28] are the elements of the transposed inverse of the matrix  $(2K_{\alpha\beta})$ , and are given by

$$G_{\alpha\beta} = i(1 - \omega^{\alpha\beta}) e^{-\pi i \alpha / N} / (2NL_\alpha),$$

while the  $K'_{\alpha\beta}$  are

$$K'_{\alpha\beta} = \sum_{\gamma=1}^{\beta} \bar{K}_{\alpha\gamma} = \frac{i e^{\pi i \alpha / N} (\omega^{-\alpha\beta} - 1)}{2 \sin(\pi \alpha / N)} L'_\alpha. \quad (15)$$

From now on take sums and products over  $\alpha, \beta, \gamma$  to be from 1 to  $N - 1$  unless otherwise indicated. Inserting some omitted primes in the relevant equation in [28], define

$$\tau_{\alpha\beta} = 2i \sum_{\gamma} G_{\gamma\alpha} K'_{\gamma\beta}. \quad (16)$$

Using (13), this gives

$$\tau_{\alpha\beta} = \frac{2i}{N} \sum_{\gamma} \frac{\cos \frac{\pi\gamma(\alpha-\beta)}{N} \sin \frac{\pi\alpha\gamma}{N} \sin \frac{\pi\beta\gamma}{N} L'_{\gamma}}{\sin \frac{\pi\gamma}{N} L_{\gamma}}. \quad (17)$$

for  $\alpha, \beta = 1, \dots, N - 1$ . The matrix with elements  $-\tau_{\alpha\beta}$  is real, symmetric and positive-definite.

## Hyperelliptic Functions

Write the variables  $u_1, \dots, u_q; v_1, \dots, v_q$  of [28] as  $w_1, \dots, w_q; s_1, \dots, s_q$ . Define new variables  $v_1, \dots, v_q$  by  $x_{\alpha} = -\cot[(v_{\alpha} + \theta)/N]$ ,  $x_1, \dots, x_q$  being as in [28]. Then for  $\alpha = 1, \dots, N - 1$ ,

$$w_{\alpha} = \sum_{\beta} \int_{\pi\beta-\pi+\theta}^{v_{\beta}} \frac{e^{i(N-2\alpha)v/N} dv}{\Delta(v)}; \quad (18)$$

$$w_{\alpha} = 2 \sum_{\beta} K_{\alpha\beta} s_{\beta}, \quad s_{\alpha} = \sum_{\beta} G_{\beta\alpha} w_{\beta}. \quad (19)$$

Writing the ordered set of variables  $\{s_1, \dots, s_{N-1}\}$  simply as  $s$ , the hyperelliptic theta function is

$$\Theta\{s\} = \sum_m \exp\{2\pi i \sum_{\alpha} m_{\alpha} s_{\alpha} + \pi i \sum_{\alpha} \sum_{\beta} m_{\alpha} \tau_{\alpha\beta} m_{\beta}\}, \quad (20)$$

the outer sum being over all values of the integers  $m = \{m_1, \dots, m_{N-1}\}$ .

More generally, let  $n = \{n_1, \dots, n_{N-1}\}$  be a given set of rational numbers, and define

$$\Theta\{s|n\} = \sum_m \exp\{2\pi i \sum_{\alpha} (m_{\alpha} + n_{\alpha}) s_{\alpha} + \pi i \sum_{\alpha} \sum_{\beta} (m_{\alpha} + n_{\alpha}) \tau_{\alpha\beta} (m_{\beta} + n_{\beta})\}. \quad (21)$$

If  $n_1, \dots, n_{N-1}$  are all integers, then obviously  $\Theta\{s|n\} = \Theta\{s\}$ . Allow one of them to be half an odd integer and define  $\Theta\{s\}_0, \dots, \Theta\{s\}_{2N-2}$  by

$$\Theta\{s\}_{2\lambda} = \Theta(s_1, \dots, s_{\lambda}, s_{\lambda+1} - \frac{1}{2}, \dots, s_{N-1} - \frac{1}{2} | 0, \dots, 0, \frac{1}{2}, 0, \dots, 0) \quad (22)$$

$$\Theta\{s\}_{2\lambda-1} = \Theta(s_1, \dots, s_{\lambda-1}, s_{\lambda} - \frac{1}{2}, \dots, s_{N-1} - \frac{1}{2} | 0, \dots, 0, \frac{1}{2}, 0, \dots, 0),$$

the arguments  $n_1, \dots, n_{N-1}$  on the RHS being zero except in position  $\lambda$ , i.e.  $n_{\lambda} = \frac{1}{2}$ , provided  $1 \leq \lambda < N$ . (For  $\Theta\{s\}_0$  all the  $n_1, \dots, n_{N-1}$  are zero.)

## The Case $N = 3$

To fix our ideas, for  $N = 3$ ,  $\Theta\{s|n\} = \Theta(s_1, s_2 | n_1, n_2)$  and these last definitions become

$$\begin{aligned}\Theta(s_1, s_2)_0 &= \Theta(s_1 - \frac{1}{2}, s_2 - \frac{1}{2} | 0, 0), & \Theta(s_1, s_2)_1 &= \Theta(s_1 - \frac{1}{2}, s_2 - \frac{1}{2} | \frac{1}{2}, 0) \\ \Theta(s_1, s_2)_2 &= \Theta(s_1, s_2 - \frac{1}{2} | \frac{1}{2}, 0), & \Theta(s_1, s_2)_3 &= \Theta(s_1, s_2 - \frac{1}{2} | 0, \frac{1}{2}) \\ \Theta(s_1, s_2)_4 &= \Theta(s_1, s_2 | 0, \frac{1}{2}).\end{aligned}$$

In this case, setting  $\varrho = \tau_{12} = \tau_{21}$ , we have  $\tau_{11} = \tau_{22} = 2\varrho$  and (20) simplifies to

$$\Theta\{s\} = \sum_{m_1, m_2} \exp\{2\pi i(m_1 s_1 + m_2 s_2) + 2\pi i\varrho(m_1^2 + m_1 m_2 + m_2^2)\}. \quad (23)$$

Changing the summation indices to  $m_1 + m_2$  and  $m_1 - m_2$  (both odd or both even), it follows that

$$\Theta\{s\} = \theta_2(\pi s_1 + \pi s_2, q^3) \theta_2(\pi s_1 - \pi s_2, q) + \theta_3(\pi s_1 + \pi s_2, q^3) \theta_3(\pi s_1 - \pi s_2, q) \quad (24)$$

where  $q = e^{2\pi i\varrho}$  and  $\theta_2(u, q)$ ,  $\theta_3(u, q)$  are the ordinary single-variable Jacobi elliptic theta functions:

$$\begin{aligned}\theta_2(u, q) &= \sum_{n=-\infty}^{\infty} e^{(2n-1)i u} q^{(2n-1)^2/4} \\ \theta_3(u, q) &= \sum_{n=-\infty}^{\infty} e^{2ni u} q^{n^2}.\end{aligned}$$

(A similar decomposition of  $\Theta\{s\}$  into sums of products of  $\theta_2, \theta_3$  functions occurs for  $N = 4$ , but not apparently for higher  $N$ .)

The case  $N = 3$  is further discussed in the Appendix.

## Identities

Once more allowing  $N$  to be an arbitrary integer greater than one, we can now write down the first two of the three identities given by Kowalevski at the end of §6 of [28], specialized to this case. They are, for all values (real or complex) of  $v_1, \dots, v_{N-1}$ ,

$$\begin{aligned}\prod_{\beta} \frac{\sin[(v_{\beta} + \theta - \pi\alpha)/N]}{\sin[(v_{\beta} - \theta + \pi)/N]} &= (-1)^{\alpha} [\Theta\{s\}_{2\alpha}/\Theta\{s\}]^2, \\ \prod_{\beta} \frac{\sin[(v_{\beta} - \theta + \pi - \pi\alpha)/N]}{\sin[(v_{\beta} - \theta + \pi)/N]} &= (-1)^{\alpha-1} [\Theta\{s\}_{2\alpha-1}/\Theta\{s\}]^2,\end{aligned} \quad (25)$$

where  $\alpha = 0, \dots, N-1$  in the first equation;  $1, \dots, N-1$  in the second.

There is a wealth of information in just these identities. For instance, taking  $v_\beta = \pi\beta - \pi + \theta$  (all  $\beta$ ), we see from (18) and (19) that  $s_\beta = 0$ , so the first identity gives

$$\frac{\sin 2\theta}{\sin[(\pi + \pi\alpha - 2\theta)/N]} = N[\Theta\{0\}_{2\alpha}/\Theta\{0\}]^2 \quad (26)$$

for  $\alpha = 0, \dots, N-1$ , writing  $\mathbf{0}$  for  $\{0, \dots, 0\}$ . Also, taking  $v_\beta = \pi\beta - \theta$ , we obtain  $w_\alpha = \sum_\beta K_{\alpha\beta}$ ,  $s_\beta = 1/2$ , and the second identity gives

$$\frac{\sin[(\pi - 2\theta)/N]}{\sin[(\pi\alpha - \pi + 2\theta)/N]} = [\Theta\{\frac{1}{2}\}_{2\alpha-1}/\Theta\{\frac{1}{2}\}]^2 \quad (27)$$

for  $\alpha = 1, \dots, N-1$ , writing  $\frac{1}{2}$  for  $\{\frac{1}{2}, \dots, \frac{1}{2}\}$ .

## Symmetries of the $\Theta$ Functions

Define

$$\varrho_\alpha = \frac{i}{N} \sum_\gamma \frac{\sin^2(\pi\alpha\gamma/N)L'_\gamma}{\sin(\pi\gamma/N)L_\gamma}. \quad (28)$$

Then  $\varrho_\alpha = \varrho_{\alpha+N} = \varrho_{N-\alpha} = \varrho_{-\alpha}$ ;  $\varrho_0 = \varrho_N = 0$ ; and from (17)

$$\tau_{\alpha\beta} = \varrho_\alpha + \varrho_\beta - \varrho_{\alpha-\beta}. \quad (29)$$

Given  $\theta$ , we can regard either (26) or (27) as defining  $\varrho_1, \dots, \varrho_{N-1}$ .

The definition (20) can be written more symmetrically by setting  $s_\alpha = r_\alpha - r_N$  ( $r_N$  arbitrary). Then

$$\Theta\{s\} = \sum_m \exp\{2\pi i \sum_{\alpha=1}^N m_\alpha r_\alpha - \pi i \sum_{\alpha=1}^N \sum_{\beta=1}^N \varrho_{\alpha-\beta} m_\alpha m_\beta\}, \quad (30)$$

the outer sum now being over all values of the  $N$  integers  $m = \{m_1, \dots, m_N\}$ , subject to the condition

$$m_1 + \dots + m_N = 0.$$

For  $\alpha$  an integer other than  $1, \dots, N-1$ , define  $s_\alpha$  by  $s_0 = s_N = 0$ ,  $s_\alpha = s_{\alpha+N} = s_{\alpha-N}$ . (This is consistent with (19).) Write  $\Theta\{s\}$  alternatively as  $\Theta\{s_\alpha\}$ , the index  $\alpha$  being understood to range over all integers, in particular over  $0, \dots, N-1$ . Define  $\delta_{\alpha\beta}$  to be one if  $\alpha = \beta \pmod{N}$ , zero otherwise. Then from (20) we readily see that  $\Theta\{s_\alpha\}$  satisfies the quasi-periodicity and evenness relations

$$\Theta\{s_\alpha + \delta_{\alpha j} - \delta_{0j}\} = \Theta\{s_\alpha\}, \quad (31a)$$

$$\Theta\{s_\alpha + \tau_{\alpha j}\} = \exp\{-2\pi i(s_j + \varrho_j)\} \Theta\{s_\alpha\} \quad (31b)$$

$$\Theta\{-s_\alpha\} = \Theta\{s_\alpha\}, \quad (31c)$$

for  $j = 0, \dots, N-1$ .

The RHS of (30) is clearly unchanged by cyclic permutations of  $r_1, \dots, r_N$  (similarly permuting  $m_1, \dots, m_N$ ), in particular by  $r_\alpha \rightarrow r_{\alpha-1}$ , so

$$\Theta\{s_\alpha\} = \Theta\{s_{\alpha-1} - s_{N-1}\}. \quad (31d)$$

Also, replacing  $m_\alpha, r_\alpha$  in (30) by  $m_{N+1-\alpha}, r_{N+1-\alpha}$ , we see that

$$\Theta\{s_\alpha\} = \Theta\{s_{N+1-\alpha} - s_1\}. \quad (31e)$$

The other functions  $\Theta\{s\}_0, \dots, \Theta\{s\}_{2N-2}$  can be expressed in terms of  $\Theta\{s\}$  and satisfy similar symmetry relations. In particular,

$$\Theta\{s_\alpha\}_1 = -i \exp\{\pi i(s_1 + \frac{1}{2}\varrho_1)\} \Theta\{s_\alpha + \frac{1}{2}(\delta_{\alpha 0} - 1 + \tau_{\alpha 1})\} \quad (32)$$

and is an odd function:

$$\Theta\{-s_\alpha\}_1 = -\Theta\{s_\alpha\}_1, \quad \Theta\{0\}_1 = 0. \quad (33)$$

## Specialization to a Single Variable $v$

Now we focus on the case which seems to be relevant to the chiral Potts model, namely when

$$\begin{aligned} v_1 &= v \\ v_\beta &= \pi\beta - \pi + \theta, \quad \beta = 2, \dots, N-1. \end{aligned} \quad (34)$$

Thus there is only one variable  $v$ , and the expression (18) simplifies to

$$w_\alpha = \int_0^v \frac{e^{i(N-2\alpha)v/N}}{\Delta(v)}, \quad (35)$$

where  $\alpha = 1, \dots, N-1$ . Obviously  $w_1, \dots, w_{N-1}$  are not independent: once one (and  $N, k$ ) is specified, the rest are determined to within a discrete set of choices. Similarly,  $s_1, \dots, s_{N-1}$  are not independent. In fact, from (25),

$$\Theta\{s\}_3 = \Theta\{s\}_5 = \dots = \Theta\{s\}_{2N-3} = 0, \quad (36)$$

which relations can be regarded as defining  $s_2, \dots, s_{N-1}$  in terms of  $s_1$ .

The function  $\Theta\{s\}_1$  does not vanish identically, but is zero when  $v = 0$ , i.e. when  $s_1 = s_2 = \dots = s_{N-1} = 0$ .

Let  $\mathcal{D}$  be the region of the complex  $v$ -plane consisting of the lower half plane, together with the real axis except for the points  $v = \pi\alpha \pm \theta$  (for all integers  $\alpha$ ). If we use the above sign conventions for  $\Delta(v)$ , and take the path of integration in (35) to be a straight line, then the  $s_\alpha$  are analytic in  $\mathcal{D}$ . So are  $u$  and  $\mu$ , but they are not uniquely defined by (3a)–(3c).

We can fix the choice of  $u$  and  $\mu$  by specifying them when  $v$  is real, between  $-\theta$  and  $\theta$ . Then they can be chosen so that  $-\pi/2 < u < \pi/2$  and  $\mu e^{-iv/N}$  is real and positive; each  $s_\alpha$  is negative imaginary, between  $-\tau_{\alpha 1}/2$  and 0.

These conventions define the  $s_\alpha, u, \mu$  uniquely for  $v \in \mathcal{D}$ . We extend these definitions by analytic continuation beyond  $\mathcal{D}$ . The points  $v = \pi\alpha \pm \theta$  are branch points, and  $s_\alpha, u, \mu$  become multi-valued functions of  $v$ . A convenient tool for handling this multi-valuedness is the following set of automorphisms.

## Automorphisms

There are five mappings (or sets of mappings) that change  $u, v, \mu$  (and hence  $x, y, a, b, c, d$ ), but leave the relations (3a)–(3d) unchanged. Including the corresponding transformation of  $s_1, \dots, s_{N-1}$ , they are

$$\begin{aligned} M_j^{(1)} &: v, u; \mu, s_\alpha \rightarrow v, u; \omega\mu, s_\alpha + \delta_{\alpha j} - \delta_{0j} \\ M_j^{(2)} &: v, u; \mu, s_\alpha \rightarrow v, u + 2j\pi; \mu, s_\alpha + \tau_{\alpha j} \\ M^{(3)} &: v, u; \mu, s_\alpha \rightarrow v, \pi - u; e^{2iv/N}\mu^{-1}, -s_\alpha \\ M^{(4)} &: v, u; \mu, s_\alpha \rightarrow v + \pi, u + \pi; \omega\mu, s_{\alpha-1} - s_{N-1} + \frac{1}{2}(\delta_{\alpha 1} + \tau_{\alpha 2} - \tau_{\alpha 1}) \\ M^{(5)} &: v, u; \mu, s_\alpha \rightarrow -v, -u; e^{-2iv/N}\mu, s_{N+1-\alpha} - s_1 - \frac{1}{2}\tau_{\alpha 1}. \end{aligned} \quad (37)$$

The mappings  $M_j^{(1)}$ ,  $M_j^{(2)}$ ,  $M^{(3)}$  manifest the fact that  $u, \mu, s_1, \dots, s_{N-1}$  are multi-valued functions of  $v$ . By deforming the path of integration in (35) to make an extra loop round the line segment  $(\pi j - \pi + \theta, \pi j - \theta)$ , we increase  $w_\alpha$  by  $2K_{\alpha j}$ , resulting in the mapping  $M_j^{(1)}$  of  $\mu$  and  $s_\alpha$ . Here  $j$  can be any integer, but only the values  $0, 1, \dots, N-1$  give distinct mappings.

Similarly, including one extra loop round each of the  $j$  line segments  $(-\theta, \theta)$ ,  $(\pi - \theta, \pi + \theta)$ ,  $(\pi j - \pi - \theta, \pi j - \pi + \theta)$  increases  $w_\alpha$  by  $2iK'_{\alpha j}$ , and results in the mapping  $M_j^{(2)}$ . The mapping  $M^{(3)}$  is obtained by negating  $A(v)$  in (35), i.e. moving onto the other Riemann sheet for  $A(v)$ .

Incrementing  $v$  by  $\pi$  replaces  $w_\alpha$  by  $\omega^{-\alpha}(w_\alpha + K_{\alpha 0} + iK'_{\alpha 1})$  and gives  $M^{(4)}$ ; negating  $v$  replaces  $w_\alpha$  by  $-w_{N-\alpha} - iK'_{\alpha 1}$  and gives  $M^{(5)}$ . Once  $v, u, \mu$  are known, then  $x, y$  and the ratios  $a : b : c : d$  are uniquely determined by (3b) and (3d). Hence we can normalize  $a, b, c, d$  so that

$$\begin{aligned} M_j^{(1)} &: x, y; a, b, c, d \rightarrow x, y; \omega a, b, c, \omega d \\ M_j^{(2)} &: \quad \quad \quad \rightarrow \omega^j x, \omega^j y; \omega^j a, \omega^j b, c, d \\ M^{(3)} &: \quad \quad \quad \rightarrow \omega/y, \omega/x; c, \omega^{1/2} d, \omega^{-1/2} a, \omega^{-1} b \\ M^{(4)} &: \quad \quad \quad \rightarrow x, \omega y, a, b, \omega^{-1} c, d \\ M^{(5)} &: \quad \quad \quad \rightarrow x^{-1}, \omega/y; d, \omega^{1/2} c, \omega^{-1/2} b, a. \end{aligned} \quad (38)$$

This makes it clear that the five mappings merely permute (and possibly negate)  $a^N, b^N, c^N, d^N$ : they leave the set of relations (1) unchanged. These automorphisms are related to the  $R, S, T, U$  of previous papers [4], [10], [20]; in particular

$$R = M^{(3)}M^{(5)}, \quad S = M^{(3)}M_1^{(2)}. \quad (39)$$

## Zeros of $a, b, c, d$

Take  $v \in \mathcal{D}$  and choose  $u, \mu$  as specified above. By applying Cauchy's theorem to the rectangle with corners  $-\theta, \pi - \theta, \pi - \theta - iR, -\theta - iR$  ( $R$  large and positive), we can verify that in the limit  $v \rightarrow -i\infty$ ,

$$w_\alpha = (\omega^\alpha K_{\alpha 1} + iK'_{\alpha 1}) / (\omega^\alpha - 1). \quad (40)$$

It follows that  $s_\alpha$  is then

$$s_\alpha = \frac{N - \alpha}{2N} + \frac{1}{2}(\varrho_{\alpha-1} - \varrho_1), \quad \alpha = 1, \dots, N, \quad (41)$$

and  $u - v = i \ln k$ .

Applying the mappings  $M^{(3)}$  and  $M^{(5)}$ , we deduce that if (for  $\alpha = 1, \dots, N$ )

$$s_\alpha = \lambda \frac{\alpha - N}{2N} + \frac{\lambda v}{2} (\varrho_{\alpha-1} - \varrho_1), \quad (42)$$

where  $\lambda, v = \pm 1$ , then  $u \rightarrow i\lambda\infty, v \rightarrow iv\infty$ . Writing the RHS of (42) as  $\zeta_{\lambda v \alpha}$ , the function  $\Theta\{s_\alpha - \zeta_{\lambda v \alpha}\}_1$  then vanishes. From (32), so therefore does

$$\Theta\{s_\alpha - \frac{1}{2}\lambda g_\alpha - \frac{1}{2}\lambda v \varrho_\alpha\},$$

where

$$\begin{aligned} g_\alpha &= \alpha/N, \quad \alpha = 0, \dots, N-1, \\ &= g_{\alpha+N}, \quad \forall \alpha. \end{aligned} \quad (43)$$

Define

$$\begin{aligned} \hat{x} &= C \Theta\{s_\alpha - \frac{1}{2}g_\alpha + \frac{1}{2}\varrho_\alpha\}/\Theta\{s_\alpha + \frac{1}{2}g_\alpha + \frac{1}{2}\varrho_\alpha\}, \\ \hat{y} &= C' \Theta\{s_\alpha - \frac{1}{2}g_\alpha - \frac{1}{2}\varrho_\alpha\}/\Theta\{s_\alpha + \frac{1}{2}g_\alpha - \frac{1}{2}\varrho_\alpha\}, \\ \hat{\mu} &= C'' e^{2\pi i \langle s \rangle} \Theta\{s_\alpha + \frac{1}{2}g_\alpha + \frac{1}{2}\varrho_\alpha\}/\Theta\{s_\alpha + \frac{1}{2}g_\alpha - \frac{1}{2}\varrho_\alpha\}, \end{aligned} \quad (44)$$

where

$$\langle s \rangle = N^{-1} \sum_{\alpha=1}^{N-1} s_\alpha$$

and  $C, C', C''$  are some constants. Then from (3b), (3c) these functions  $\hat{x}, \hat{y}, \hat{\mu}$  have the same zeros and poles at  $v = \pm i\infty$  as do  $x, y, \mu$ , respectively.

Further, using (31a) and (31b), we find that  $\hat{x}, \hat{y}, \hat{\mu}$  transform under  $M_j^{(1)}$  and  $M_j^{(2)}$  in the same way as  $x, y, \mu$ ; the ratios  $\hat{x}/x, \hat{y}/y, \hat{\mu}/\mu$  being invariant under these automorphisms. It follows at once that the ratios are only two-valued functions of  $v$ , being uniquely determined by  $v$  and  $A(v)$ . They can be regarded as defined on two Riemann sheets of the complex  $v$ -plane.

By itself, (31d) does not appear sufficient to establish the invariance of the ratios under  $M^{(4)}$ : one presumably needs to also use the relations (36). However, the  $s_\alpha$  are invariant under  $M^{(4)}_j$ , so the ratios are periodic functions of  $v$ , of period  $N\pi$ .

The following statements have not yet been rigorously proved, but (guided by the small- $k'$  limit) they appear to be correct: (i) The ratios  $\hat{x}/x, \hat{y}/y, \hat{\mu}/\mu$  remain finite and non-zero as  $v \rightarrow \pm i\infty$  on either sheet. (ii) The  $\Theta$  functions in (44) are non-zero for all finite  $v$ .

Assuming these statements to be true, it follows that  $\hat{x}/x$  is bounded on both Riemann sheets; analytic except possibly at the branch points, where it is finite. By a straightforward extension of Liouville's theorem in complex variable theory,

it follows that  $\hat{x}/x$  is a constant. We can choose  $C$  so that this ratio is one. Arguing similarly for  $y$  and  $\mu$ , it follows that

$$x = \hat{x}, \quad y = \hat{y}, \quad \mu = \hat{\mu}. \quad (45)$$

From the mapping  $M^{(3)}$ , the constants  $C, C', C''$  must satisfy

$$CC' = \omega, \quad C' = e^{i\pi/N} CC''^2. \quad (46)$$

Also, when  $v = \theta$  and  $s_1, \dots, s_{N-1} = 0$ ,  $x\mu = CC' = e^{i\pi/2N}$ . It follows that there exists a single constant  $\chi$  such that

$$\begin{aligned} C &= e^{i\pi/N} \chi^{-1}, \quad C' = e^{i\pi/N} \chi, \\ C'' &= e^{-i\pi/2N} \chi, \end{aligned} \quad (47)$$

and

$$\begin{aligned} a:b:c:d &= e^{i\pi/2N} e^{2\pi i \langle s \rangle} \Theta \left\{ s_\alpha - \frac{1}{2}g_\alpha + \frac{1}{2}\varrho_\alpha \right\} : \chi e^{i\pi/N} \Theta \left\{ s_\alpha - \frac{1}{2}g_\alpha - \frac{1}{2}\varrho_\alpha \right\} : \\ &\quad \Theta \left\{ s_\alpha + \frac{1}{2}g_\alpha - \frac{1}{2}\varrho_\alpha \right\} : \chi e^{-i\pi/2N} e^{2\pi i \langle s \rangle} \Theta \left\{ s_\alpha + \frac{1}{2}g_\alpha + \frac{1}{2}\varrho_\alpha \right\}. \end{aligned} \quad (48)$$

Thus, in this normalization, and to within factors that are constant or exponential in  $s_1, \dots, s_{N-1}$ , the variables  $a, b, c, d$  are hyperelliptic theta functions. (For the case  $N = 2$ , which corresponds to the Ising model, they are proportional to the four Jacobi theta functions  $\theta_1, \theta_4, \theta_3, \theta_2$ , all with the same argument.)

### Expressions for $k, k', \chi$

We can evaluate  $\chi$  by taking  $s_\alpha = \pm \frac{1}{2}(\varrho_{\alpha-1} - \varrho_1) \pm \frac{1}{2}(N - \alpha)/N$ , in which case one of  $a, b, c, d$  vanishes. Substituting the forms (48) into (1), we obtain three distinct relations:

$$(k'/k)^{1/N} = i e^{-\pi i (\varrho_1 + \dots + \varrho_{N-1})/N} \Theta \{g_\alpha\}_1 / \Theta \{\varrho_\alpha\}_1, \quad (49)$$

$$\begin{aligned} k^{1/N} &= \chi \Theta \{\varrho_\alpha\}_1 / \Theta \{\varrho_\alpha - g_\alpha\}_1, \\ &= \chi^{-1} \Theta \{\varrho_\alpha\}_1 / \Theta \{\varrho_\alpha + g_\alpha\}_1. \end{aligned} \quad (50)$$

We can regard these equations as defining  $\chi$ , and providing alternative expressions for  $k, k'$ . In fact, for  $N = 2, 3$  and  $4$  we have observed that 13-digit numerical calculations are fitted by (48)–(50), with

$$\chi = i e^{i(2-N)\theta/N}, \quad (51)$$

As explained in the Appendix, we have also verified this formula to 201 terms in a series expansion for  $N = 3$ . It seems likely that it is exactly correct, for all  $N$ . Remembering that  $g_\alpha$  is real and  $\varrho_\alpha$  is pure imaginary, (51) is equivalent to

$$\arg \Theta \{g_\alpha + \varrho_\alpha\}_1 = (N - 2)\theta/N \quad (52)$$

an intriguingly simple conjecture.

## Appendix: The case $N = 3$

Here we consider the case  $N = 3$ , when the hyperelliptic  $\Theta\{s\} = \Theta(s_1, s_2)$  functions can be expressed, using (24), in terms of ordinary single-variable Jacobi theta functions of nomes  $q$  and  $q^3$ , where  $q = e^{2\pi i \varrho}$ .

Setting  $\phi = \pi/2 - \theta$ , the relation (26), for  $\alpha = 0$ , becomes

$$4 \sin^2(2\phi/3) = \frac{3[\Theta(0,0)^2 - \Theta(\frac{1}{2}, \frac{1}{2})^2]}{\Theta(0,0)^2}, \quad (\text{A1})$$

and hence, using (2),

$$k'^2 = 1 - k^2 = \sin^2 \phi = 27q\{1 - 15q + 171q^2 + \dots\}. \quad (\text{A2})$$

Here we regard these equations as defining  $\phi, k, k'$  as functions of  $q$ , for  $|q| < 1$ .

From these definitions, we should like to directly verify the above results (49), (50) and the conjecture (51). These can be written, using (32), as

$$(k'/k)^{1/3} = -i e^{\pi i(l+\varrho)/3} \Theta(-\frac{1}{6} + \varrho, \frac{1}{6} + \frac{1}{2}\varrho) / \Theta(\frac{1}{2}, \frac{1}{2} + \frac{1}{2}\varrho), \quad (\text{A3})$$

$$\begin{aligned} k^{1/3} &= \chi e^{\pi i/3} \Theta(\frac{1}{2}, \frac{1}{2} + \frac{1}{2}\varrho) / \Theta(\frac{1}{6}, -\frac{1}{6} - \frac{1}{2}\varrho), \\ &= \chi^{-1} e^{-\pi i/3} \Theta(\frac{1}{2}, \frac{1}{2} + \frac{1}{2}\varrho) / \Theta(\frac{1}{6}, -\frac{1}{6} + \frac{1}{2}\varrho), \end{aligned} \quad (\text{A4})$$

$$\chi = e^{i(\pi+\phi)/3}. \quad (\text{A5})$$

Define

$$Q(q) = \prod_{n=1}^{\infty} (1 - q^n) = \sum_{n=-\infty}^{\infty} (-1)^n q^{n(3n-1)/2}, \quad (\text{A6})$$

and write  $Q(q^r)$  simply as  $Q_r$ . Then from (24) we can prove that

$$\begin{aligned} \Theta(0,0) + \Theta(\frac{1}{2}, \frac{1}{2}) &= 2(Q_2 Q_6)^5 / (Q_1 Q_3 Q_4 Q_{12})^2, \\ \Theta(0,0) - \Theta(\frac{1}{2}, \frac{1}{2}) &= 8q (Q_4 Q_{12})^2 / (Q_2 Q_6), \\ \Theta(\frac{1}{2}, \frac{1}{2} + \frac{1}{2}\varrho) &= \Theta(\frac{1}{2}, \frac{1}{2} - \frac{1}{2}\varrho) = Q(q)^2. \end{aligned} \quad (\text{A7})$$

We also conjecture the following five identities

$$\Theta(\frac{1}{2}, \frac{1}{2}) = (Q_1 Q_3)^2 / (Q_2 Q_6), \quad (\text{A8a})$$

$$\Theta(0,0)^3 Q_1^3 Q_3^3 = Q_1^{12} + 27q Q_3^{12}, \quad (\text{A8b})$$

$$\Theta(-\frac{1}{6} + \varrho, \frac{1}{6} + \frac{1}{2}\varrho) = (1 - e^{-2\pi i/3}) Q_3^2, \quad (\text{A8c})$$

$$\Theta(\frac{1}{6}, -\frac{1}{6} \pm \frac{1}{2}\varrho) = [Q_2^3 Q_3 / (Q_1 Q_6)] \pm i(3q)^{1/2} [Q_1 Q_6^3 / (Q_2 Q_3)], \quad (\text{A8d})$$

$$\Theta(\frac{1}{6}, -\frac{1}{6} + \frac{1}{2}\varrho) \Theta(\frac{1}{6}, -\frac{1}{6} - \frac{1}{2}\varrho) = \Theta(0,0) Q_1 Q_3. \quad (\text{A8e})$$

We have verified these conjectures to order  $q^{200}$  in series expansions in powers of  $q$ , using Fortran with integer arithmetic. The hyperelliptic  $\Theta$  functions

herein, as well as  $Q(q)$ , have  $q$ -expansions that are sparse, with the few non-zero coefficients being of order one. By writing the identities solely in terms of sums of products of such functions (with *no* divisions), we avoided integer overflows: for instance, to this order the largest coefficient in the expansion of  $Q(q)^{12}$  is only 5,187,456.

The desired formulae (A3), (A4), (A5) follow from (A1), (A2) and (A6)–(A8e), so we have verified them to order  $q^{200}$ . The equation (A3) becomes very simple:

$$(k'/k)^{1/3} = 3^{1/2} q^{1/6} \{Q(q^3)/Q(q)\}^2. \quad (\text{A9})$$

### Note

A number of colleagues have pointed out that the conjectures (A8a)–(A8b) can be proved, either by the general theory of modular functions, or from specific identities already in the literature. In particular, the author is indebted to G.E. Andrews for showing their connection with the theory of generalized Frobenius partitions [29].

### References

1. H. Au-Yang, B.M. McCoy, J.H.H. Perk, S.Tang, M.L. Yan: Phys. Lett. A **123** (1987) 219–223
2. B.M. McCoy, J.H.H. Perk, S. Tang, C.H. Sah: Phys. Lett. A **125** (1987) 9–14
3. H. Au-Yang, B.M. McCoy, J.H.H. Perk, S. Tang: In: M. Kashiwara and T. Kawai (eds.) Algebraic Analysis vol. 1. Academic Press, New York 1988, pp. 29–39
4. R.J. Baxter, J.H.H. Perk, H. Au-Yang: Phys. Lett. A **128** (1988) 138–142
5. R.J. Baxter: J. Stat. Phys. **52** (1988) 639–667
6. R.J. Baxter: Phys. Lett. A **133** (1988) 185–189
7. J.H.H. Perk Proc. of Symposia in Pure Mathematics, vol. **49**, part 1 (1989) 341–354
8. G. Albertini, B.M. McCoy, J.H.H. Perk, S. Tang: Nucl. Phys. B **314** (1989) 741–763
9. G. Albertini, B.M. McCoy, J.H.H. Perk: Adv. Stud. Pure Math. **19** (1989) 1–55
10. H. Au-Yang, J.H.H. Perk: Adv. Stud. Pure Math. **19** (1989) 57–94
11. R.J. Baxter: Adv. Stud. Pure Math. **19** (1989) 95–116
12. G. Albertini, B.M. McCoy, J.H.H. Perk: Phys. Lett. A **135** (1989) 159–166
13. G. Albertini, B.M. McCoy, J.H.H. Perk: Phys. Lett. A **139** (1989) 204–212
14. R.J. Baxter: Phys. Lett. A **140** (1989) 155–157
15. R.J. Baxter: J. Stat. Phys. **57** (1989) 1–39
16. V.B. Mateev and A.O. Smirnov: Some comments on the solvable chiral Potts model. Preprint 1989
17. H. Itoyama, B.M. McCoy, J.H.H. Perk: Intnl. J. Mod. Phys. B **4** (1990) 995–1001
18. D. Bernard and V. Pasquier: Intnl. J. Mod. Phys. B **4** (1990) 913–927
19. V.V. Bazhanov, Yu. G. Stroganov: J. Stat. Phys. **59** (1990) 799–817
20. R.J. Baxter, V.V. Bazhanov, J.H.H. Perk: Intnl. J. Mod. Phys. B **4** (1990) 803–870
21. G. Albertini, B.M. McCoy: Nucl. Phys. B **350** (1991) 745–788
22. R.J. Baxter: Phys. Lett. A **146** (1990) 110–114
23. B.M. McCoy, S. Roan: Phys. Lett. A **150** (1990) 347–354
24. R.J. Baxter: Calculation of the eigenvalues of the transfer matrix of the chiral Potts model. For the Proceedings of the Fourth Asia Pacific Conference 1990 (World Scientific)

25. V.V. Bazhanov, R.M. Kashaev: Cyclic L-operator related with a 3-state R-matrix, preprint RIMS-702 (1990). To appear in Comm. Math. Phys.
26. V.V. Bazhanov, R.M. Kashaev, V.V. Mangazeev, Yu.G. Strogarov:  $(Z_N)^{n-1}$  generalization of the chiral Potts model. To appear in Comm. Math. Phys.
27. R.J. Baxter: J. Stat. Phys. **63** (1991) 509–529
28. S. Kowalevski: Acta Mathematica **12** (1889) 177–232
29. G.E. Andrews: Memoirs of the American Math. Soc. **49** (1984) no. 301, iv.



# Abstract Compact Group Duals, Operator Algebras and Quantum Field Theory

Sergio Doplicher

Dipartimento di Matematica, Università di Roma “La Sapienza”, I-00185 Rome, Italy

## 1. Introduction

In this report we will review joint work with J. E. Roberts on the existence of a unique compact group of internal symmetries, and of field operators, on which the group acts as global gauge transformations, which commute or anticommute in the normal way at spacelike separations [16].

These results are deduced from first principles (essentially, the postulate of locality) involving solely the observable quantities, for local Quantum Theories without massless particles on the four dimensional Minkowski space.

On the mathematical side, these results stem out of a new duality theory for compact groups which has been developed for this purpose. While the classical theory of Tannaka and Krein (cf. e.g. [22]) characterizes the dual of a compact group within the category of finite dimensional vector spaces, the new theory characterizes such a dual within abstract categories.

More specifically, our abstract compact group duals will be strict symmetric monoidal  $C^*$ -categories with conjugates, having subobjects and direct sums, where the selfintertwiners of the monoidal unit reduce to  $\mathbb{C}$  [15].

Each such category can be “locally” realized as a full subcategory of  $\text{End } \mathfrak{A}$  (cf. Sect. 2), i.e. as an action on a  $C^*$ -algebra.

Now if a category  $\mathcal{T}$  with the mentioned properties acts as a full subcategory of  $\text{End } \mathfrak{A}$  on a  $C^*$ -algebra  $\mathfrak{A}$  with centre  $\mathbb{C} \cdot I$ , a key construction [14] provides a cross product  $\mathfrak{A} \times \mathcal{T}$  containing  $\mathfrak{A}$  as a  $C^*$ -subalgebra with trivial relative commutant, and a compact group  $G$  of automorphisms of  $\mathfrak{A} \times \mathcal{T}$  having  $\mathfrak{A}$  as the fixed point subalgebra, such that the objects of  $\mathcal{T}$  are the restrictions to  $\mathfrak{A}$  of inner endomorphisms of  $\mathfrak{A} \times \mathcal{T}$ , induced by Hilbert spaces with support  $I$  in  $\mathfrak{A} \times \mathcal{T}$ ; these Hilbert spaces together with  $\mathfrak{A}$  generate  $\mathfrak{A} \times \mathcal{T}$  and are  $G$ -spaces.

Assigning to each object of  $\mathcal{T}$  the representation of  $G$  on the corresponding Hilbert space in  $\mathfrak{A} \times \mathcal{T}$  and to each arrow the  $G$ -module map given by multiplication in  $\mathfrak{A} \times \mathcal{T}$  defines a functor of  $\mathcal{T}$  into the category of  $G$ -Hilbert spaces in  $\mathfrak{A} \times \mathcal{T}$  which is an isomorphism of strict symmetric monoidal categories of  $\mathcal{T}$  with a category of continuous unitary finite dimensional representations of  $G$ .

These central results, when applied “locally” to an abstract category  $\mathcal{T}$  in the earlier mentioned class, yield at the same time the existence of a unique compact

group  $G$  and of an isomorphism of  $\mathcal{T}$  with a representation category of  $G$  as strict symmetric monoidal categories [15, Theorem 6.1].

By completely different methods, a parallel result has been independently established by Pierre Deligne in a recent paper [7] characterizing as abstract categories the categories of representations of an algebraic group on finite dimensional vector spaces over a field of characteristic zero.

In the Quantum Field Theory context, we let  $\mathfrak{U}$  be the  $C^*$ -algebra generated by all local observables (Sect. 3). The superselection structure associated to the given vacuum state  $\omega_0$  on  $\mathfrak{U}$  is determined by  $\mathfrak{U}$  itself and can be described by a category  $\mathcal{T}$  fulfilling exactly the axioms mentioned above. In the important case of “localizable charges”,  $\mathcal{T}$  is actually a full subcategory of  $\text{End } \mathfrak{U}$ . The monoidal structure induced from  $\text{End } \mathfrak{U}$  is related to composition of superselection charges; the symmetry to the intrinsic notion of statistics of superselection sectors; the existence of conjugates is a consequence of particle-antiparticle symmetry of superselection quantum numbers. All these pieces of structure have been deduced, essentially from the locality principle, in earlier works [9, 10], and can be summarized saying that superselection structure is described by a full subcategory  $\mathcal{T}$  of  $\text{End } \mathfrak{U}$  with the above properties. The corresponding cross product  $\mathfrak{U} \times \mathcal{T}$  describes *field operators* and the dual action  $G$  on  $\mathfrak{U} \times \mathcal{T}$  provides the compact group of internal symmetries. With some technical modifications needed to take care of “quantum topological charges” this construction yields the more general results described in Section 3.

## 2. Abstract Compact Group Duals and Operator Algebras

Which *abstract* categories can appear as the dual object of a compact group? The general answer can be summarized in the following theorem [15]:

**2.1. Theorem.** *Let  $\mathcal{T}$  be a strict symmetric monoidal  $C^*$ -category with conjugates, having subobjects and direct sums, such that the selfinterwiners of the monoidal unit are  $\mathbb{C}$ . There is a unique compact group  $G$  and an isomorphism of symmetric monoidal  $C^*$ -categories of  $\mathcal{T}$  with a category of finite dimensional continuous unitary representations of  $G$ .*

Rather than giving the formal definitions involved in this statement, it is easier to see how they arise naturally in an important model.

Let  $\mathfrak{U}$  be a  $C^*$ -algebra with centre  $\mathbb{C} \cdot I$ . We will denote by  $\text{End } \mathfrak{U}$  the category whose objects are the unital endomorphisms of  $\mathfrak{U}$  and whose arrows  $(\varrho, \varrho')$  for object  $\varrho, \varrho'$  are their intertwiners in  $\mathfrak{U}$ :

$$(\varrho, \varrho') = \{T \in \mathfrak{U} \mid T\varrho(A) = \varrho'(A)T, A \in \mathfrak{U}\}. \quad (1)$$

$\text{End } \mathfrak{U}$  is a  $C^*$ -category with the linear structure, norm and  $*$  induced by  $\mathfrak{U}$  and composition of arrows  $T \circ S$ , when defined, given by product in  $\mathfrak{U}$ . Further, it is a monoidal  $C^*$ -category, where the monoidal operations are composition in the

semigroup with identity  $\iota$  (the identity automorphism) of unital endomorphisms, and on arrows

$$\begin{aligned} T \in (\varrho, \varrho'), S \in (\sigma, \sigma') &\rightarrow T \times S \in (\varrho\sigma, \varrho'\sigma'); \\ T \times S = T\varrho(S) &= \varrho'(S)T. \end{aligned} \quad (2)$$

These operations are always defined, *strictly associative*, and fulfill obvious compatibility relations with composition, so that  $\text{End } \mathfrak{A}$  is a *strict* monoidal  $C^*$ -category. The selfinterwiners of the monoidal unit  $\iota$  are the elements of the centre of  $\mathfrak{A}$ , i.e.  $(\iota, \iota) = \mathbb{C} \cdot I$ .

We are interested in *full monoidal subcategories*  $\mathcal{T}$  of  $\text{End } \mathfrak{A}$  (i.e.  $\mathcal{T}$  is specified by a semigroup of unital endomorphisms with identity  $\iota$  and all arrows given by (1)) which are the model for a representation category of a group. The monoidal operations on  $\mathcal{T}$  should behave like the *tensor product* in a subcategory of Hilbert spaces where  $\otimes$  is *strictly associative*. This suggests the following

**2.2. Definition.** A symmetry  $\varepsilon$  for a strict monoidal  $C^*$ -category  $\mathcal{T}$  is an assignment of a unitary  $\varepsilon(\varrho, \sigma) \in (\varrho\sigma, \sigma\varrho)$  for each pair of objects  $\varrho, \sigma$ , so that if  $R \in (\varrho, \sigma)$ ,  $R' \in (\varrho', \sigma')$ ,

$$\varepsilon(\sigma, \sigma') \circ R \times R' = R' \times R \circ \varepsilon(\varrho, \varrho') \quad (3)$$

and

$$\varepsilon(\sigma, \varrho) \circ \varepsilon(\varrho, \sigma) = I_{\varrho\sigma} \quad (4)$$

$$\varepsilon(\iota, \varrho) = \varepsilon(\varrho, \iota) = I_\varrho \quad (5)$$

$$\varepsilon(\varrho\sigma, \tau) = \varepsilon(\varrho, \tau) \times I_\sigma \circ I_\varrho \times \varepsilon(\sigma, \tau). \quad (6)$$

The symmetry  $\varepsilon$  defines canonical unitary representations  $\varepsilon_q^{(n)}$  of the permutation groups  $\mathbb{P}(n)$  of  $n$  objects with values in  $(\varrho^n, \varrho^n)$ ,  $n = 2, 3, \dots$ , obtained assigning  $I_{\varrho^{r-1}} \times \varepsilon(\varrho, \varrho) \times I_{\varrho^{n-r-1}}$  to the exchange of  $r$  with  $r+1$  (where  $I_\sigma \in (\sigma, \sigma)$  is the identity interwiner).

Crucial for *compactness* is the existence of conjugates. The following definition is modelled on the properties of complex conjugates of finite dimensional group representations.

**2.3. Definition.** The strict symmetric monoidal  $C^*$ -category  $(\mathcal{T}, \varepsilon)$  has conjugates if to each object  $\varrho$  of  $\mathcal{T}$  there is an object  $\bar{\varrho}$  and arrows  $R, \bar{R}$  fulfilling

$$R \in (\iota, \bar{\varrho}\varrho), \bar{R} \equiv \varepsilon(\bar{\varrho}, \varrho) \circ R, \quad (7)$$

$$\bar{R}^* \times I_\varrho \circ I_\varrho \times R = I_\varrho, \quad (8)$$

$$R^* \times I_{\bar{\varrho}} \circ I_{\bar{\varrho}} \times \bar{R} = I_{\bar{\varrho}}. \quad (9)$$

Eventually we will say that  $\mathcal{T}$  has *subobjects* if each non zero selfadjoint projection  $E \in (\varrho, \varrho)$  is the range of an isometry  $W \in (\sigma, \varrho)$  for some object  $\sigma$ ;  $\mathcal{T}$  has *direct sums* if any pair of objects  $\sigma, \sigma'$  are subobjects of some object  $\varrho$ , associated to complementary projections in  $(\varrho, \varrho)$ .

Now the statement of Theorem 2.1 is explained. When  $\mathcal{T}$  is actually a full subcategory of  $\text{End } \mathfrak{U}$ , and  $\mathfrak{U}$  has centre  $\mathbb{C} \cdot I$ , slightly more general axioms suffice, replacing existence of conjugates with abundance of objects with “determinant one”. We can define determinants in an abstract category  $\mathcal{T}$  as above along the following lines.

It can be shown [15, Sect. 2] that each object  $\varrho$  has an *integer dimension*  $d(\varrho)$  given by

$$R^* \circ R = d(\varrho) \cdot I \quad (10)$$

where  $R$  obeys the conjugate equations (7,8,9). The objects with dimension one are shown to form a *group* and their equivalence classes form an *abelian group*  $\mathcal{G}_0$ . The *determinant map*

$$\text{Objects } \mathcal{T} \rightarrow \mathcal{G}_0$$

assigns to  $\varrho$  the class of the subobject (automatically of dimension one) of  $\varrho^{d(\varrho)}$  associated to the projection  $\varepsilon_{\varrho}^{d(\varrho)}(A_{d(\varrho)}) \in (\varrho^{d(\varrho)}, \varrho^{d(\varrho)})$ , where  $A_{d(\varrho)}$  is the totally antisymmetric projection in the group algebra of  $\mathbb{P}(d(\varrho))$  [15, Sect. 3].

The determinant map has the properties

$$\det \varrho_1 \oplus \varrho_2 = \det \varrho_1 \cdot \det \varrho_2 \quad (11)$$

$$\det \bar{\varrho} = (\det \varrho)^{-1} \quad (12)$$

$$\det \varrho_1 \varrho_2 = 1 \quad \text{if } \det \varrho_1 = \det \varrho_2 = 1. \quad (13)$$

Objects with determinant one are called *special*. An object  $\varrho$  of dimension  $d$  is special if and only if there is an *isometry*  $R \in (\iota, \varrho^d)$  such that

$$R \circ R^* = \varepsilon_{\varrho}^{(d)}(A_d) \quad (14)$$

$$R^* \times I_{\varrho} \circ I_{\varrho} \times R = (-1)^{d-1} \frac{1}{d} I_{\varrho}.$$

We conclude that  $(\mathcal{T}, \varepsilon)$  is *specially directed* in the sense that each finite set of objects  $\varrho_1, \dots, \varrho_n$  are dominated by a special object  $\varrho$ : it suffices to choose

$$\varrho = \varrho_1 \oplus \dots \oplus \varrho_n \oplus \bar{\varrho}_n \oplus \dots \oplus \bar{\varrho}_1, \quad (15)$$

which will be special by (11) and (12).

Now we can state precisely in which sense  $\mathcal{T}$  is locally described by a full subcategory of  $\text{End } \mathfrak{U}$  [15, Sect. 4].

**2.4. Theorem.** *Let  $\mathcal{T}$  be as in Theorem 2.1,  $\varrho$  a special object and  $\mathcal{T}_{\varrho}$  the full subcategory with objects  $\{\iota, \varrho, \varrho^2, \dots\}$ . There is a  $C^*$ -algebra  $\mathcal{O}_{\varrho}$  with centre  $\mathbb{C} \cdot I$ , a strict symmetric monoidal full subcategory  $\mathcal{T}_{\hat{\varrho}}$  of  $\text{End } \mathcal{O}_{\varrho}$  with objects the powers of a special object  $\hat{\varrho}$ , and a functor of  $\mathcal{T}_{\varrho}$  onto  $\mathcal{T}_{\hat{\varrho}}$  which is an isomorphism of symmetric monoidal categories.*

To construct  $\mathcal{O}_{\varrho}$  one considers, for each  $k \in \mathbb{Z}$ , the inductive limit  $\mathcal{O}_{\varrho}^k$  of the Banach spaces  $(\varrho^r, \varrho^{r+k})$ ,  $r, r+k \in \mathbb{N}_0$ , under the maps

$$T \in (\varrho^r, \varrho^{r+k}) \rightarrow T \times I_{\varrho} \in (\varrho^{r+1}, \varrho^{r+k+1}).$$

The composition of maps and the adjoint in  $\mathcal{T}$  make  $\mathcal{O}^k$ ,  $k \in \mathbb{Z}$ , into a  $\mathbb{Z}$ -graded  $*\text{-algebra}$  which can be completed in a unique way into a  $C^*$ -algebra  $\mathcal{O}_\varrho$  such that the grading is given by a continuous action of  $\mathbb{T}$ .

Denoting by  $i: T \rightarrow i(T)$  the embedding of  $(\varrho^*, \varrho^s)$  into  $\mathcal{O}_\varrho$  we have

$$\begin{aligned} i(T \times I_\varrho) &= i(T), \\ i(I_\varrho \times T) &= \hat{\varrho} \circ i(T), \end{aligned} \tag{16}$$

which defines the endomorphism  $\hat{\varrho}$ .

Theorem 2.4 brings the abstract duality problem into its natural context, non commutative operator algebras.

The main result, crucial both for the proof of Theorem 2.1 and for the application to QFT, deals with a monoidal full subcategory  $\mathcal{T}$  of  $\text{End } \mathfrak{A}$  (hence  $\mathcal{T}$  is specified by its objects, a semigroup  $\Lambda$  of unital endomorphisms with identity  $i$ ) which has a symmetry  $\varepsilon$  (but  $\mathcal{T}$  does not need to have subobjects and direct sums: if it did,  $\mathcal{T}$  would have conjugates if and only if it were specially directed, cf. the comments above and [12, 9, 10]).

Such a category arises naturally if  $\mathfrak{A}$  is embedded in a  $C^*$ -algebra  $\mathfrak{B}$  as the fixed points under a compact group  $G$  of automorphisms of  $\mathfrak{B}$  and  $\mathfrak{A}' \cap \mathfrak{B} = \mathbb{C} \cdot I$ . In this case a finite dimensional *Hilbert space*  $H$  in  $\mathfrak{B}$  (i.e.  $\psi^* \psi' \in \mathbb{C} \cdot I$  if  $\psi, \psi' \in H$ ) with support  $I$  (i.e. the left annihilator of  $H$  is zero) induces an endomorphism  $\varrho_H$  of  $\mathfrak{A}$  if and only if  $H$  is  $G$  stable; then

$$H = \{\psi \in \mathfrak{B} \mid \psi A = \varrho_H(A)\psi, \quad A \in \mathfrak{A}\}. \tag{17}$$

Then  $H$  carries a unitary representation  $U_H$  of  $G$  given by  $U_H(g)\psi = g(\psi)$ ,  $\psi \in H$ . The representations so obtained and their intertwiners form a category denoted  $\mathcal{U}(\mathfrak{B}, G)$ . The linear operators  $(H, H')$  between Hilbert spaces  $H, H'$  as above form a linear subspace of  $\mathfrak{B}$  (spanned by  $H'H^*$ ) and their  $G$ -invariant elements  $(H, H')_G$  lie in  $\mathfrak{A}$  and are the arrows of  $\mathcal{U}(\mathfrak{B}, G)$ . Also  $(H, H')_G \subset (\varrho_H, \varrho_{H'})$  and  $U_H \rightarrow \varrho_H$ ,  $T \in (U_H, U_{H'}) \rightarrow T \in (\varrho_H, \varrho_{H'})$  sets up an isomorphism of monoidal categories of  $\mathcal{U}(\mathfrak{B}, G)$  with the full subcategory  $\mathcal{S}(\mathfrak{B}, G)$  of  $\text{End } \mathfrak{A}$  with objects  $\{\varrho_H, H \text{ a finite dimensional } G\text{-Hilbert space in } \mathfrak{B} \text{ with support } I\}$ . The monoidal structure in  $\mathcal{U}(\mathfrak{B}, G)$  is induced by the operator product in  $\mathfrak{B}$  which defines a strictly associative tensor product of Hilbert spaces in  $\mathfrak{B}$ .

This monoidal structure has a natural symmetry  $\theta$  defined by

$$\begin{aligned} \theta_{H,H'} &\in (HH', H'H), \\ \theta_{H,H'} \psi \psi' &= \psi' \psi; \quad \psi \in H, \psi' \in H', \end{aligned} \tag{18}$$

where actually  $\theta_{H,H'} \in \mathfrak{A}$ . Then  $(\mathcal{U}(\mathfrak{B}, G), \theta)$  and  $(\mathcal{S}(\mathfrak{B}, G), \theta)$  are isomorphic as strict symmetric monoidal categories [13]. Representations with determinant one correspond to special objects.

Every category of finite dimensional continuous representations of a compact group  $G$  can be embedded in a category  $\mathcal{U}(\mathfrak{B}, G)$  as above [13, Sect. 7].

**2.5. Definition.** We will say that  $(\mathcal{T}, \varepsilon)$  is a symmetric monoidal subcategory of  $(\mathcal{S}(\mathfrak{B}, G), \theta)$  if  $\mathcal{T}$  is a monoidal subcategory of  $\mathcal{S}(\mathfrak{B}, G)$  and  $\varepsilon$  is just the restriction

of the symmetry  $\theta$  to objects in  $\mathcal{T}$ . In this case we will denote by  $H(\mathcal{T})$  the Hilbert spaces in  $\mathfrak{B}$  inducing the objects of  $\mathcal{T}$ .

The main result of [14] provides a converse construction given  $\mathfrak{A}$  and  $\mathcal{T}$ .

**2.6. Theorem.** *Let  $\mathfrak{A}$  be a  $C^*$ -algebra with centre  $\mathbb{C} \cdot I$  and  $(\mathcal{T}, \varepsilon)$  a symmetric, specially directed full monoidal subcategory of  $\text{End } \mathfrak{A}$ .*

*There is a  $C^*$ -algebra  $\mathfrak{B}$  and a strongly compact group  $G$  of automorphisms of  $\mathfrak{B}$  such that*

- 1)  $\mathfrak{A} = \mathfrak{B}^G$ ;
- 2)  $\mathfrak{A}' \cap \mathfrak{B} = \mathbb{C} \cdot I$ ;
- 3)  $(\mathcal{T}, \varepsilon)$  is a symmetric monoidal subcategory of  $(\mathcal{S}(\mathfrak{B}, G), \theta)$ ;
- 4)  $\mathfrak{B}$  is generated as a  $C^*$ -algebra by  $\mathfrak{A}$  and  $H(\mathcal{T})$ .

*As a consequence,  $G$  is the stabilizer of  $\mathfrak{A}$  in  $\text{Aut } \mathfrak{B}$  and  $(\mathfrak{B}, G)$  is unique up to isomorphisms that leave  $\mathfrak{A}$  pointwise fixed.*

The  $C^*$ -algebra  $\mathfrak{B}$  may be called the crossed product  $\mathfrak{A} \times \mathcal{T}$  and  $G$  the dual action. The group  $G$  would be automatically compact provided we replace it by its strong closure if needed [13].

Note that the Theorem provides at the same time  $G$  and a representation category of  $G$  isomorphic to  $\mathcal{T}$ , namely the image in  $\mathcal{U}(\mathfrak{B}, G)$  of the subcategory  $\mathcal{T}$  of  $\mathcal{S}(\mathfrak{B}, G)$ .

The special case  $\mathcal{T} = \mathcal{T}_\varrho$ ,  $\varrho$  a special object yields a compact Lie group. Combined with Theorem 2.4 this yields a compact Lie group  $G_\varrho$  for each special object of a category  $\mathcal{T}$  as in Theorem 2.1; the group  $G$  of Theorem 2.1 is obtained as the projective limit of the compact Lie groups  $G_\varrho$ .

The proof of Theorem 2.6 is reduced to the case  $\mathcal{T} = \mathcal{T}_\varrho$ . Indeed, if  $\hat{\mathcal{T}}_\varrho$  is the full subcategory of  $\text{End } \mathfrak{A}$  with objects the objects of  $\mathcal{T}$  dominated by a special object  $\varrho$ , then  $\mathfrak{A} \times \hat{\mathcal{T}}_\varrho$  is easily constructed from  $\mathfrak{A} \times \mathcal{T}_\varrho$  and  $\mathfrak{A} \times \mathcal{T}$  can be obtained as the inductive limit of the  $C^*$ -algebras  $\mathfrak{A} \times \hat{\mathcal{T}}_\varrho$ .

The construction of  $\mathfrak{A} \times \mathcal{T}_\varrho$  is then central, and we outline the main ideas involved. Being generated by a special  $G$ -Hilbert module  $H$  and by  $\mathfrak{A}$ ,  $\mathfrak{A} \times \mathcal{T}_\varrho$  will then contain both  $\mathfrak{A}$  and the Cuntz algebra [6]  $\mathcal{O}_d$  generated by  $H$ , where  $d = d(\varrho) = \dim H$ , with intersection the fixed points  $\mathcal{O}_G$  in  $\mathcal{O}_d$  under the action of  $G$ . This action is canonical in the sense that it is induced by the unitary representation of  $G$  on  $H$ ,  $U(g)\psi = g(\psi)$ ,  $\psi \in H$ ,  $g \in G$ . Thus, if the system  $(\mathfrak{B}, G)$  is given, we have monomorphisms  $\mu : \mathcal{O}_G \rightarrow \mathfrak{A}$ ,  $\pi : \mathfrak{A} \rightarrow \mathfrak{B}$ ,  $\zeta : \mathcal{O}_d \rightarrow \mathfrak{B}$  such that

$$\zeta(\psi)\pi(A) = \pi\varrho(A)\zeta(\psi), \quad A \in \mathfrak{A}, \quad \psi \in H \tag{19}$$

$$\mu \circ \sigma = \varrho \circ \mu \tag{20}$$

(where  $\sigma$  is the endomorphism of  $\mathcal{O}_d$  induced by the generating Hilbert space  $H$ ) such that the following diagram is commutative

$$\begin{array}{ccc}
 \mathfrak{A} & \xrightarrow{\pi} & \mathfrak{B} \\
 \uparrow \mu & & \uparrow \varsigma \\
 \mathcal{O}_G & \longrightarrow & \mathcal{O}_d .
 \end{array} \tag{21}$$

As a consequence of  $\mathfrak{A}' \cap \mathfrak{B} = \mathbb{C} \cdot I$ , we have

$$\mu((\sigma^r, \sigma^s)) = (\varrho^r, \varrho^s) ; r, s \in \mathbb{N}_0 . \tag{22}$$

But  $G$  and  $\mu$  are *not* known a priori. The solution of the problem involves three steps.

a) Since  $\varrho$  is special, identifying  $G$  with its faithful representations by  $d \times d$  unitary matrices given by the action on  $H$ , we have  $G \subset SU(d)$ . Then  $\mathcal{O}_{SU(d)} \subset \mathcal{O}_G$ . Now  $\mu_0 \equiv \mu|_{\mathcal{O}_{SU(d)}}$  is actually *determined* by the data,  $\varepsilon_\varrho^{(n)}, n = 2, 3, \dots$  and  $R$  fulfilling (14) for  $\varrho$ .

b) Assuming that  $G$  and  $\mu$  are given fulfilling (20) and  $\mu((\sigma^r, \sigma^s)) \subset (\varrho^r, \varrho^s) ; r, s \in \mathbb{N}_0$  (cf. (22)), find a *universal solution* to the diagram (21) with condition (19).

c) Apply this construction to  $SU(d), \mu_0$ . We get a  $C^*$ -algebra  $\tilde{\mathfrak{B}}$  with an action of  $SU(d)$  with  $\mathfrak{A}$  as fixed points but, since (22) does not hold for  $\mu_0$  in general, we will not have  $\mathfrak{A}' \cap \tilde{\mathfrak{B}} = \mathbb{C} \cdot I$ .

However  $\mathfrak{A}' \cap \tilde{\mathfrak{B}}$  is *commutative*, and the system  $(\tilde{\mathfrak{B}}, SU(d))$  can be obtained as the induced  $C^*$  system from a unique system  $(\mathfrak{B}, G)$  where now  $\mathfrak{A} = \mathfrak{B}^G$  fulfills  $\mathfrak{A}' \cap \mathfrak{B} = \mathbb{C} \cdot I$ .

Then  $G$  and  $\mu$  fulfilling (20), (22) are obtained and we can identify  $(\mathfrak{B}, G)$  with the universal solution of step b).

Step a) involves the  $C^*$ -algebraic version of the theory of invariants for  $SU(d)$  [12].

If  $S \in H^d$  is a normalized totally antisymmetric vector,  $S \in \mathcal{O}_{SU(d)}$  and, together with  $\theta^{(n)}(p), p \in \mathbb{P}(n), n = 2, 3, \dots$  (where  $\theta^{(n)}(p) \in (H^n, H^n)$  permutes the factors in  $H^n$  by  $p$ ), generates  $\mathcal{O}_{SU(d)}$  as a  $C^*$ -algebra, and this  $C^*$ -algebra is *simple*.

Now if  $\varrho$  is special symmetric,  $\varepsilon_\varrho^{(n)}$  and  $R$  (cf. eq. (14)) will generate a  $C^*$ -algebra isomorphic to  $\mathcal{O}_{SU(d)}$ , and we can define a unique monomorphism  $\mu_0$  of  $\mathcal{O}_{SU(d)}$  into  $\mathfrak{A}$ , fulfilling (20), by

$$\begin{aligned}
 \mu_0(\theta^{(n)}(p)) &= \varepsilon_\varrho^{(n)}(p), \quad p \in \mathbb{P}(n), \quad n = 2, 3, \dots \\
 \mu_0(S) &= R .
 \end{aligned}$$

Step b) involves first an algebraic construction [14, Sect. 2]. With  ${}^\circ\mathcal{O}_d$  the dense  $*$  subalgebra of  $\mathcal{O}_d$  generated by  $H$ , and  ${}^\circ\mathcal{O}_G$  its  $G$  invariant part, we can regard  ${}^\circ\mathcal{O}_d$  and  $\mathfrak{A}$  as  ${}^\circ\mathcal{O}_G$ -bimodules letting  $X \in {}^\circ\mathcal{O}_G$  act by multiplication in  ${}^\circ\mathcal{O}_d$  and by multiplication with  $\mu(X)$  in  $\mathfrak{A}$ . The  ${}^\circ\mathcal{O}_G$ -module tensor product  $\mathfrak{A} \otimes_{{}^\circ\mathcal{O}_d} {}^\circ\mathcal{O}_d$  can be made into an algebra by requiring that  $I \otimes_{\mu_0} H$  induces  $\varrho$  on  $\mathfrak{A} \otimes I$ , i.e. setting

$$\begin{aligned}
 A \otimes_\mu X \cdot B \otimes_\mu Y &= A \varrho^n(B) \otimes_\mu XY, \\
 A, B \in \mathfrak{A}, \quad X \in H^n, \quad Y \in {}^\circ\mathcal{O}_d .
 \end{aligned} \tag{23}$$

To define the  $*$ -operation it is essential to use the fact that  $S \in {}^0\mathcal{O}_G$  since  $G \subset SU(d)$ . Since  $\psi^* = S^*S\psi^* = S^*j(\psi)$ , where  $j$  is the antilinear map of  $H$  into  $H^{d-1}$ ,  $j(\psi) = d(-1)^{d-1}\psi^*S$  (as  $S^*\sigma(S) = (-1)^{d-1}\frac{1}{d}I$ , cf. also (14)), we can define the adjoint map on the generators  $\mathfrak{A} \otimes_{\mu} H$  by

$$\begin{aligned}(A \otimes_{\mu} \psi)^* &= (A \otimes_{\mu} I \cdot I \otimes_{\mu} \psi)^* = I \otimes_{\mu} \psi^* \cdot A^* \otimes_{\mu} I = \\ &= I \otimes_{\mu} S^* j(\psi) \cdot A^* \otimes_{\mu} I = \mu(S^*) \varrho^{d-1}(A^*) \otimes_{\mu} j(\psi).\end{aligned}\quad (24)$$

Along these lines it is possible to make  $\mathfrak{A} \otimes_{\mu} {}^0\mathcal{O}_d$  into a  $*$ -algebra and there is a unique  $C^*$ -norm extending the norm of  $\mathfrak{A}$  which is continuous for the action of  $G$ . The completion  $\mathfrak{A} \otimes_{\mu} \mathcal{O}_d$  in this norm is the desired universal solution to (19), (21) [14, Sect. 3].

Taking now  $\mu = \mu_0$  given by step a) in this construction we find a  $C^*$ -system  $(\tilde{\mathfrak{B}}, SU(d))$  with fixed points  $\mathfrak{A}$  such that, as a consequence of the conditions  $\mu_0(\theta(p)) = \varepsilon_q(p)$  we have [13, Theorem 5.2]

$$\mathfrak{A}' \cap \tilde{\mathfrak{B}} = \tilde{\mathfrak{B}}' \cap \tilde{\mathfrak{B}}. \quad (25)$$

The desired system  $(\mathfrak{B}, G)$  is now constructed in step c) noting that by (25) and  $\mathfrak{A}' \cap \mathfrak{A} = \mathbb{C} \cdot I$ ,  $SU(d)$  acts ergodically on  $\tilde{\mathfrak{B}}' \cap \tilde{\mathfrak{B}}$ , hence, by compactness, transitively on its spectrum. Picking  $\phi$  in that spectrum,  $G$  will be the stabilizer of  $\phi$ ; the smallest closed two sided ideal  $J$  in  $\tilde{\mathfrak{B}}$  containing  $\ker \phi$  will be  $G$  stable hence  $G$  will act on  $\mathfrak{B} = \tilde{\mathfrak{B}}/J$  providing the desired solution [11].

Now  $G$  is known and the identification of  $\mathcal{O}_G$  with  $\mathfrak{A} \cap \mathcal{O}_d$  in  $\mathfrak{B}$  defines  $\mu$ , fulfilling (22) since  $\mathfrak{A}' \cap \mathfrak{B}$  is now  $\mathbb{C} \cdot I$ . The corresponding universal solution of (19), (21)  $(\mathfrak{A} \otimes_{\mu} \mathcal{O}_d, G)$  is now isomorphic to  $(\mathfrak{B}, G)$  [14].

Theorem 2.6 admits spatial versions of direct use in QFT. If  $m$  denotes the conditional expectation  $\mathfrak{A} \times \mathcal{T} \rightarrow \mathfrak{A}$  given by integration of the action of  $G$  over the normalized Haar measure, a faithful representation  $\pi_0$  of  $\mathfrak{A}$  induces a representation  $\pi$  of  $\mathfrak{B}$  via  $m$  [14, Sect. 6].

If  $\pi_0(\mathfrak{A})$  is a factor  $M$ ,  $N \equiv \pi(\mathfrak{B})''$  will be a factor with  $M' \cap N = \mathbb{C} \cdot I$ ; conversely if  $N$  is a factor with separable predual and  $G \subset \text{Aut } N$  is compact for the strong topology of the action on the predual, such that  $M = N^G$  is infinite with trivial relative commutant, then  $N$  arises from  $M$  as above [14, Theorem 7.1].

We could generalize  $\text{End } \mathfrak{A}$  to  $\text{Bimod } \mathfrak{A}$ , a strict monoidal  $C^*$ -category whose objects are homomorphisms  $\varrho$  of  $\mathfrak{A}$  into square matrix algebras over  $\mathfrak{A}$  and where an arrow in  $(\varrho, \sigma)$  is a matrix over  $\mathfrak{A}$  intertwining  $\varrho$  and  $\sigma$  [15, Sect. 1].

$\text{Bimod } \mathfrak{A}$  can be locally identified with  $\text{End } \tilde{\mathfrak{A}}$  for a suitable  $C^*$ -algebra  $\tilde{\mathfrak{A}}$ . If the object  $\varrho$  is a unital mapping of  $\mathfrak{A}$  into  $M_n(\mathfrak{A})$ , let  $\tilde{\mathfrak{A}}$  be the  $C^*$ -tensor product of  $\mathfrak{A}$  with the UHF algebra  $M_n(\mathbb{C}) \otimes M_n(\mathbb{C}) \otimes \dots$ , and define an endomorphism  $\tilde{\varrho}$  of  $\tilde{\mathfrak{A}}$  by

$$\begin{aligned}\tilde{\varrho}(A \otimes B_1 \otimes B_2 \otimes B_3 \otimes \dots) &= \varrho(A) \otimes B_1 \otimes B_2 \otimes \dots, \\ B_i &\in M_n(\mathbb{C}), \quad A \in \mathfrak{A}.\end{aligned}\quad (26)$$

It can be proved that  $\mathcal{T}_{\varrho}$  is isomorphic to  $\mathcal{T}_{\tilde{\varrho}}$  [15, Sect. 5]. Actions described by  $\text{Bimod } \mathfrak{A}$  arise, e.g. if the  $C^*$ -system  $(\mathfrak{B}, G)$  does not have  $G$ -Hilbert spaces. A unitary matrix  $X \in \mathfrak{B} \otimes M_n(\mathbb{C})$  such that

$$g \otimes \iota(X) = X \cdot I \otimes v(g), \quad g \in G, \quad (27)$$

where  $v$  is a continuous unitary  $n$ -dimensional representation of  $G$ , defines an object of  $\text{Bimod } \mathfrak{A}$ ,  $\mathfrak{A} = \mathfrak{B}^G$ , by

$$\varrho(A) = XA \otimes IX^*, \quad A \in \mathfrak{A}. \quad (28)$$

By the foregoing device this action could be equally well described by an action in  $\text{End } \mathfrak{A}$ .

On the other side, every category of continuous unitary finite dimensional representations of a compact group  $G$  can be embedded in  $\text{End } \mathfrak{A}$  for a *simple*  $C^*$ -algebra  $\mathfrak{A}$ , which is obtained as the fixed point subalgebra under a canonical action of  $G$  on an infinite tensor product of Cuntz algebras [13, Sect. 7].

As a variant of that construction one can consider the fixed points  $\mathcal{O}_{U(G)}$  in the Cuntz algebra  $\mathcal{O}_\infty$  under the action  $\alpha$  of a group  $G$  defined by a unitary representation  $U$  of  $G$  on the infinite dimensional generating Hilbert space  $H$ , by

$$\alpha_g(\psi) = U(g)\psi; \quad \psi \in H, \quad g \in G.$$

It can be shown [5] that this action is *ergodic* if and only if  $U$  admits no non zero finite dimensional subrepresentation, and is *prime* if  $U(G)$  has compact strong closure. If furthermore  $U$  is the direct sum of finite dimensional subrepresentations with determinant one,  $\mathcal{O}_{U(G)}$  is *simple* [5], generalizing to  $\mathcal{O}_d$  the analogous result for  $\mathcal{O}_d$  [12].

By taking, for each metrizable compact group  $G$ , the representation  $U$  to be the doubled left regular representation  $\lambda$ , we get a *simple*  $C^*$ -algebra  $\mathcal{O}_{\lambda(G)}$ , and each continuous isomorphism  $G_1 \rightarrow G_2$  induces a canonical isomorphism  $\mathcal{O}_{\lambda(G_2)} \rightarrow \mathcal{O}_{\lambda(G_1)}$ .

Current research deals with the question: when are the  $\mathcal{O}_{\lambda(G)}$  non isomorphic?

### 3. Operator Algebras and Quantum Field Theory

Quantum Mechanics says that observables generate a  $C^*$ -algebra  $\mathfrak{A}$  [31]; local quantum theory, providing a general basis to Quantum Field Theory, requires that  $\mathfrak{A}$  is generated by *local observables* [21]. Local observables are specified by an *inclusion preserving* map  $\mathcal{O} \rightarrow \mathfrak{A}(\mathcal{O})$  from the set  $\mathcal{K}$  of double cones (bounded set obtained intersecting past with future open light cones) in Minkowski space to von Neumann algebras acting on a separable Hilbert space  $\mathcal{H}_0$ . This map is extended to any  $S \subset \mathbb{R}^4$  defining  $\mathfrak{A}(S)$  as the  $C^*$ -algebra generated by  $\mathfrak{A}(\mathcal{O})$ ,  $\mathcal{O} \subset S$ ,  $\mathcal{O} \in \mathcal{K}$ ; so that  $\mathfrak{A} = \mathfrak{A}(\mathbb{R}^4)$ .

The  $C^*$ -algebra  $\mathfrak{A}$  is assumed to be *irreducible*, i.e. the defining representation  $\pi_0$  of  $\mathfrak{A}$  describes a single *superselection sector*, which should contain the *vacuum state*  $\omega_0$  induced by a unit vector  $\Omega \in \mathcal{H}_0$ .

Other superselection sectors will be described by other unitary equivalence classes of irreducible representations of  $\mathfrak{A}$ .

As a manifestation of Einstein causality, if the double cones  $\mathcal{O}_1, \mathcal{O}_2$  cannot be joined by any signal travelling at most with the speed of light, i.e.  $\mathcal{O}_1$  is included

in the spacelike complement  $\mathcal{O}'_2$  of  $\mathcal{O}_2$ , observables localized in  $\mathcal{O}_1$  and  $\mathcal{O}_2$  should not interfere with one another. Quantum Mechanics says that they will commute; hence, for each  $\mathcal{O} \in \mathcal{K}$ , the *locality principle* says that

$$\mathfrak{A}(\mathcal{O}) \subset \mathfrak{A}(\mathcal{O}')'. \quad (29)$$

In QFT, fields are basically associated with points or, possibly, with loops or infinite strings. While a double cone is an appropriate neighbourhood of a point, an appropriate neighbourhood of a string is the cone joining a point to a double cone at spacelike infinity. The set of such *spacelike cones* is denoted by  $\mathcal{J}$ .

Locality is often sharpened to *duality*

$$\mathfrak{A}(\mathcal{S})^- = \mathfrak{A}(\mathcal{S}')' \quad (30)$$

where  $\mathcal{S}$  is a double cone (so that weak closure on the left-hand side is unnecessary) or a spacelike cone.

Duality is related to the absence of spontaneously broken gauge symmetries [29]. Otherwise, *essential duality* [30], requiring that the net  $\mathcal{O} \in \mathcal{K} \rightarrow \mathfrak{A}^d(\mathcal{O}) \equiv \mathfrak{A}(\mathcal{O})'$  fulfills (30) for each double cone, would hold quite generally, as a consequence of a result of Bisognano and Wichmann [1], whenever there are underlying Wightman fields.

Our last assumption is Property B: if  $\mathcal{S}, \tilde{\mathcal{S}}$  are either double cones or spacelike cones such that  $\mathcal{S} + \mathcal{N} \subset \tilde{\mathcal{S}}$  for some neighborhood  $\mathcal{N}$  of the origin, non zero self adjoint projections in  $\mathfrak{A}(\mathcal{S})^-$  are equivalent to  $I \bmod \mathfrak{A}(\tilde{\mathcal{S}})^-$ , i.e.

$$\mathcal{O} \neq E = E^*E \in \mathfrak{A}(\mathcal{S})^- \Rightarrow E = WW^*, W^*W = I, W \in \mathfrak{A}(\tilde{\mathcal{S}})^-. \quad (31)$$

If we would assume the basic requirements of translation covariance and spectrum condition in the vacuum sector, Property B would be a consequence, as discovered by Borchers [2].

These few axioms on the inclusion preserving map  $\mathcal{O} \rightarrow \mathfrak{A}(\mathcal{O})$  (irreducibility, duality and Property B), though too general on their own to characterize physically the vacuum sector, turn out to be sufficient for a discussion of superselection structure tailored to theories without massless particles.

In such a theory one particle excitations of the vacuum would be associated to (factorial) representations  $\pi$  of  $\mathfrak{A}$  which are translationally covariant, with energy momentum spectrum (the spectrum of the representation of the translation group) in the forward light cone starting with an isolated hyperboloid of positive mass (massive particle representations). If  $\pi$  describes an excitation of the vacuum  $\omega_0$ ,  $\omega_0$  can be obtained as a mean of  $\pi$  over translations.

However, Buchholz and Fredenhagen [4] derived, from the above assumptions, a much stronger relation expressing that  $\pi$  can be localized in spacelike cones:

$$\pi|_{\mathfrak{A}(\mathcal{C})} \cong \pi_0|_{\mathfrak{A}(\mathcal{C})}, \mathcal{C} \in \mathcal{J}. \quad (32)$$

The general definition and classification of *statistics* [9, 10] (cf. also J. E. Roberts contribution to [26], and comments below) can be extended to representations fulfilling (32), ([4], cf. also [16, Sect. 4]) and moreover massive particle

representations have necessarily *finite statistics* [17] (so that in particular  $\pi$  is of type I).

These basic facts motivate, in our general frame where covariance and spectrum condition are not assumed, the following

**3.1. Definition.** *In a theory specified by the inclusion preserving irreducible map  $\mathcal{O} \rightarrow \mathfrak{A}(\mathcal{O})$  fulfilling duality and Property B, the superselection sectors are the unitary equivalence classes of irreducible representations  $\pi$  of  $\mathfrak{A}$  fulfilling (32) with finite statistics.*

Superselection sectors are usually associated with charged field operators which connect the vacuum sector to the other sectors. These fields might be attached to points (Wightman fields) or to strings from points to spacelike infinity, in presence of “quantum topological charges”. We can formalize these notions in the following definition [16].

**3.2. Definition.** *An extended field system with gauge symmetry  $(\pi, G, \mathfrak{F})$  consists of*

- a representation  $\pi$  of  $\mathfrak{A}$  on a Hilbert space  $\mathcal{H}$  containing  $\pi_0$  as a subrepresentation on  $\mathcal{H}_0 \subset \mathcal{H}$ ;
- a strongly compact group  $G$  of unitaries on  $\mathcal{H}$  leaving  $\mathcal{H}_0$  pointwise fixed;
- an inclusion preserving map  $\mathfrak{F}$  from spacelike cones to von Neumann algebras on  $\mathcal{H}$  such that for each  $\mathcal{C} \in \mathcal{J}$ , the  $g \in G$  induce automorphisms  $\alpha_g$  of  $\mathfrak{F}(\mathcal{C})$ , with fixed points  $\pi(\mathfrak{A}(\mathcal{C}))^-$ ; moreover  $\mathcal{H}_0$  is cyclic for  $\mathfrak{F}(\mathcal{C})$  and the union of  $\mathfrak{F}(\mathcal{C} + a)$  for all  $a \in \mathbb{R}^4$  is irreducible;  $\mathfrak{F}(\mathcal{C})$  commutes with  $\pi(\mathfrak{A}(\mathcal{O}))$  if  $\mathcal{O}$  is a double cone spacelike to  $\mathcal{C}$ .

The system  $(\pi, G, \mathfrak{F})$  is normal if the map  $\mathfrak{F}$  obeys graded local commutativity for the  $\mathbb{Z}_2$ -grading defined by a central element  $k \in G$  with square the identity.

The system is complete if every irreducible representation  $\tilde{\pi}$  of  $\mathfrak{A}$  fulfilling (32) with finite statistics is a subrepresentation of  $\pi$ , i.e.  $\pi$  describes all superselection sectors.

The main result of [16] states that such a system can be canonically constructed from the net  $\mathcal{O} \rightarrow \mathfrak{A}(\mathcal{O})$  of local observables.

**3.3. Theorem.** *Let the irreducible, inclusion preserving map  $\mathcal{O} \rightarrow \mathfrak{A}(\mathcal{O})$  fulfill duality and Property B; there exists a complete, normal extended field system with gauge symmetry, and this system is unique up to unitary equivalence.*

Unitary equivalence of two systems  $(\pi_i, G_i, \mathfrak{F}_i)$  acting on  $\mathcal{H}_i$ ;  $i = 1, 2$ , is expressed by a unitary operator  $W$  of  $\mathcal{H}_1$ , onto  $\mathcal{H}_2$  intertwining  $\pi_1$  and  $\pi_2$ , the maps  $\mathfrak{F}_1$  and  $\mathfrak{F}_2$ , and taking  $G_1$  onto  $G_2$ :

$$\begin{aligned} W\pi_1(A) &= \pi_2(A)W, \quad A \in \mathfrak{A}; \\ W\mathfrak{F}_1(\mathcal{C}) &= \mathfrak{F}_2(\mathcal{C})W, \quad \mathcal{C} \in \mathcal{J}; \\ WG_1 &= G_2W. \end{aligned}$$

An important subclass of superselection sectors, carrying “localizable charges”, is obtained by requiring (32) to hold for *double cones*:

$$\pi_\xi|_{\mathfrak{A}(\mathcal{O})} \cong \pi_0|_{\mathfrak{A}(\mathcal{O})}, \quad \mathcal{O} \in \mathcal{K}. \quad (33)$$

This subclass is described by a unique complete normal *local field system* with gauge symmetry  $(\pi_l, G_l, \mathfrak{F}_l)$ , where now  $\mathfrak{F}_l$  is an inclusion preserving map from *double cones* to von Neumann algebras on a Hilbert space  $\mathcal{H}_l$ , with irreducible range and each  $\mathfrak{F}_l(\mathcal{O})$  is cyclic on  $\mathcal{H}_0$ .

There is a closed normal subgroup  $N$  of  $G$  identifying  $\mathcal{H}_l$  with the  $N$ -fixed vectors in  $\mathcal{H}$  and  $\pi_l$  with  $\pi|_{\mathcal{H}_l}$ ; the von Neumann algebra generated by the  $\mathfrak{F}_l(\mathcal{O})$  with  $\mathcal{O} \subset \mathcal{C}$  is identified with the restriction to  $\mathcal{H}_l$  of the  $N$ -fixed points in  $\mathfrak{F}(\mathcal{C})$ ,  $\mathcal{C} \in \mathcal{J}$  [16].

Superselection sectors described by representations fulfilling (32) but not (33) are often said to carry “quantum topological charges”.

Superselection sectors are in 1–1 correspondence with the spectrum  $\hat{G}$  of  $G$ . Quantum topological charges are carried by sectors corresponding to representations which are non trivial on  $N$ .

The superselection structure of localizable charges is described by a full subcategory  $\mathcal{T}$  of  $\text{End } \mathfrak{A}$  [9, 10].

The objects of  $\mathcal{T}$  are *localized morphisms* of  $\mathfrak{A}$  (i.e. they act trivially in the spacelike complement of some double cone which, by choosing an equivalent morphism, can be arbitrarily placed) with *finite statistics* in the following sense.

The full subcategory of  $\text{End } \mathfrak{A}$  with objects the localized morphisms has a *unique symmetry*  $\varepsilon$  which is  $I$  for spacelike separated morphisms [9]. The canonically associated representations  $\varepsilon_\varrho^{(n)}$  of  $\mathbb{P}(n)$  describe the *statistics* of  $\varrho$ : they permute the factors in the *product state vectors* in the representation  $\varrho^n$  which carry the “charges” of  $\varrho$  localized in  $n$  mutually spacelike double cones, whilst the induced *product states* are totally symmetric.

If  $\varrho$  is irreducible, the irreducible representations of  $\mathbb{P}(n)$  occurring in  $\varepsilon_\varrho^{(n)}$  are precisely those with at most  $d(\varrho) \in \mathbb{N}$  antisymmetrizations (resp. symmetrizations) and  $\varrho$  is *paraBose* (resp. *paraFermi*) of order  $d(\varrho)$ , or everyone and  $\varrho$  has infinite statistics ( $d(\varrho) = \infty$ ) [9].

Every representation  $\pi_\xi$  of  $\mathfrak{A}$  fulfilling (33) is unitarily equivalent to a localized morphism composed with  $\pi_0$ . In view of applications to massive theories on the four-dimensional spacetime we can disregard infinite statistics ([17]; see also [10; Appendix]).

We can change  $\varepsilon$  to a “Bosonized” symmetry  $\hat{\varepsilon}$  (differing from  $\varepsilon$  on irreducible elements only by sign) so that  $(\mathcal{T}, \hat{\varepsilon})$  is a strict symmetric monoidal  $C^*$ -category with *conjugates* [10, 15]. The complete normal field system with gauge symmetry can be constructed from the covariant representation of the  $C^*$ -system  $(\mathfrak{A} \times \mathcal{T}, G)$  induced from  $\pi_0$  via  $m$  [16, Sect. 3].

If only essential duality holds, we can construct the complete normal local field system with gauge symmetry from the dual net  $\mathfrak{A}^d$ . The stabilizer of  $\pi(\mathfrak{A}^d)$  in  $\text{Aut } \mathfrak{F}$  is  $G$ , the unbroken part of the gauge group. The full gauge group, possibly including *broken symmetries*, can be defined as the stabilizer  $\mathcal{G}$  of  $\pi(\mathfrak{A})$  in  $\text{Aut } \mathfrak{F}$  [16]. It can be shown that each  $\gamma \in \mathcal{G}$  leaves the local field algebras  $\mathfrak{F}(\emptyset)$  globally stable and  $\mathcal{G}$  can be in part analyzed in terms of degeneracy of the vacuum, arising from different extensions of  $\omega_0$  to  $\mathfrak{A}^d$  [3].

The duality theory reviewed in Sect. 2 suggests a natural generalization. If we drop the equation (4) in Definition 2.2 of a symmetry and we symmetrize the condition (6) adding

$$\varepsilon(\varrho, \sigma\tau) = I_\sigma \times \varepsilon(\varrho, \tau) \circ \varepsilon(\varrho, \sigma) \times I_\tau \quad (6')$$

we obtain a *strict braided* monoidal category [25].

**Problem** [15, 16]. *Which class of mathematical objects has as abstract duals the strict braided monoidal  $C^*$ -categories with conjugates, with subobjects and direct sums, such that  $(1, 1) = \mathbb{C}$ ? If such a category  $\mathcal{T}$  acts on a  $C^*$ -algebra  $\mathfrak{A}$  with centre  $\mathbb{C} \cdot I$  as a full subcategory of  $\text{End } \mathfrak{A}$ , can we generalize Theorem 2.6 to a construction of a cross product  $\mathfrak{A} \times \mathcal{T}$ ?*

This natural route to define a kind of “quantum groups” is required in order to extend the previous results to theories on low dimensional spacetimes.

To describe statistics in this case we have to replace the permutation groups by braid groups [18, 19, 20]. The category  $\mathcal{T}$  describing superselection structure is strict monoidal with a unitary braiding. Solution to the problem above would clarify which “quantum groups” may appear as internal symmetries of low dimensional models [28] (cf. also most contributions to [26]).

There are remarkable relations between the analysis of braid statistics of two dimensional theories and the Jones theory of invariant polynomials for knots and of index of  $W^*$ -inclusions [23, 24]. These relations have been discovered by R. Longo who, extending the notion of Jones index to inclusions of type III factors, showed that [27]

$$d(\varrho)^2 = \text{ind}(\varrho(\mathfrak{A}(\emptyset)), \mathfrak{A}(\emptyset)) \quad (34)$$

when  $\varrho$  is localized in  $\emptyset$ ; in two dimensional theories  $d(\varrho)$  need not be an integer but its values are limited by the equation (39). The same equation shows that  $d(\varrho)$  is related to a quantum version of the Fredholm index of an endomorphism (cf. also the comments in [8]).

## References

1. Bisognano, J.J., Wichmann, E.H.: On the duality condition for quantum fields. *J. Math. Phys.* **17** (1976) 303–321
2. Borchers, H.J.: A Remark on a theorem of B. Misra. *Comm. Math. Phys.* **4** (1967) 315–323
3. Buchholz, D., Doplicher, S., Longo, R., Roberts, J.E.: Broken symmetries and degeneracy of the vacuum in quantum field theory. Preprint in preparation

4. Buchholz, D., Fredenhagen, K.: Locality and the structure of particle states. *Comm. Math. Phys.* **84** (1982) 1–54
5. Ceccherini, T., Pinzari, C.: Canonical actions on  $\mathcal{O}_\infty$ . *J. Funct. Anal.* (to appear)
6. Cuntz, J.: Simple  $C^*$ -algebras generated by isometries. *Comm. Math. Phys.* **57** (1977) 173–185
7. Deligne, P.: Categories Tannakiennes. *Grothendieck Festschrift*. Birkhäuser, Basel 1990 (to appear)
8. Doplicher, S.: Problems in quantum field theory and operator algebras. In: Araki, H., Moore, C.C., Stratila, S., Voiculescu, D. (eds.) *Operator Algebras and their connections with Ergodic Theory and Topology*. (Lecture Notes in Mathematics, vol. 1132). Springer, Berlin Heidelberg New York 1985
9. Doplicher, S., Haag, R., Roberts, J.E.: Local observable and particle statistics I. *Comm. Math. Phys.* **23** (1971) 199–230
10. Doplicher, S., Haag, R., Roberts, J.E.: Local observables and particle statistics II. *Comm. Math. Phys.* **35** (1974) 49–85
11. Doplicher, S., Roberts, J.E.: A remark on compact automorphism groups of  $C^*$ -algebras. *J. Funct. Anal.* **66** (1986) 67–72
12. Doplicher, S., Roberts, J.E.: Duals of compact Lie groups realized in the Cuntz algebras and their actions on  $C^*$ -algebras. *J. Funct. Anal.* **74** (1987) 96–120
13. Doplicher, S., Roberts, J.E.: Compact group actions on  $C^*$ -algebras. *J. Operator Theory* **19** (1988) 283–305
14. Doplicher, S., Roberts, J.E.: Endomorphism of  $C^*$ -algebras, cross products and duality for compact groups. *Ann. Math.* **130** (1989) 75–119
15. Doplicher, S., Roberts, J.E.: A new duality theory for compact groups. *Invent. math.* **89** (1989) 157–218
16. Doplicher, S., Roberts, J.E.: Why there is a field algebra with a compact gauge group describing the superselection structure in particle physics. *Comm. Math. Phys.* **131** (1990) 51–107
17. Fredenhagen, K.: On the existence of antiparticles. *Comm. Math. Phys.* **79** (1981) 141–151
18. Fredenhagen, K., Rehren, K.H., Schroer, B.: Superselection sectors with braid group statistics and exchange algebras. *Comm. Math. Phys.* **125** (1989) 201–226
19. Fröhlich, J., in Non perturbative quantum field theory. G. 'tHooft et al. (eds.) Plenum Press, New York 1988
20. Fröhlich, J., Marchetti, P.A.: Quantum field theory of vortices and anyons. *Comm. Math. Phys.* **121** (1989) 177–224
21. Haag, R., Kastler, D.: An algebraic approach to quantum field theory. *J. Math. Phys.* **5** (1964) 848–861
22. Hewitt, E., Ross, K.A.: Abstract harmonic analysis II. Springer, Berlin Heidelberg New York 1970
23. Jones, V.F.R.: Index for subfactors. *Invent. math.* **72** (1983) 1–25
24. Jones, V.F.R.: Hecke algebra representations of Braid groups and link polynomials. *Ann. Math.* **126** (1987) 335–388
25. Joyal, A., Street, R.: Braided monoidal categories. *McQuarie Mathematics Reports No. 860081* (1986)
26. Kastler, D. (ed.): Algebraic theory of superselection sectors. World Science Publ. Co., Singapore 1990
27. Longo, R.: Index of subfactors and statistics of quantum fields I. *Comm. Math. Phys.* **126** (1989) 217–247; II. *Comm. Math. Phys.* **130** (1990) 285–310
28. Mack, G., Schomerus, V.: Conformal field algebras with quantum symmetry from the theory of superselection sectors. *Comm. Math. Phys.* **134** (1990) 139

29. Roberts, J.E.: Spontaneously broken gauge symmetries and superselection rules. *Proceedings of the International School of Mathematical Physics, Camerino 1974*, Gallavotti, G. (ed.). Università di Camerino (1976)
30. Roberts, J.E.: Net cohomology and its applications to field theory. In: L. Streit (ed.) *Quantum Fields – Algebras, Processes*. Springer, Berlin Heidelberg New York 1980
31. Segal, I.E.: Postulates for general quantum mechanics. *Ann. Math.* **48** (2) (1947) 930–948



# Introduction to Constructive Quantum Field Theory

*Joel Feldman*

University of British Columbia, Vancouver, BC, Canada, V6T 1Y4

## 1. Introduction

Quantum field theories (QFT) are used to model physical systems that share two common features. Firstly, they are of atomic or subatomic scale, so that they exhibit quantum mechanical behaviour. Secondly, they involve fields. To a physicist a field is an observable (i.e. something you measure) that is a function on space-time. Two examples are the electric and magnetic fields. Indeed the best known QFT is Quantum Electrodynamics (QED), which models the interaction of electrons with the electromagnetic field.

In addition to the obvious fields, it is often wise to associate fields with particles. For example, at the high energies typical of QED electrons and positrons are continually being created and destroyed. So, rather than attempt to explicitly keep track of how many electrons there are at all times, one invents a field, called the electron field from which one may calculate the charge and current densities. The former tells you the positions of all electrons and the latter tells you their velocities.

An electron field is also used in modelling another physical system that will appear later in this talk. The system is a crystal consisting of a large population of electrons interacting with each other and with a lattice of stationary or almost stationary ions. As there is a fixed number of electrons per ion and as the lattice is effectively of infinite size the population is really huge. Fortunately, almost all of these electrons do nothing of interest. They just sit there forming what is called the Fermi sea. However, there are small numbers of electrons continually jumping out of and back into the Fermi sea. Once again an electron field is introduced to keep track of them all.

I will use the Gårding-Wightman axioms [SW] to tell you what a quantum field theory is from a mathematical point of view. These axioms are really designed for relativistic quantum field theories so some of them are not appropriate for our current purposes and I will just mention those in passing.

The first axiom sets up the state space.

1. *State Space.* There exists a separable Hilbert space  $F$ .

The vectors (actually rays) of  $F$  are the states of our system. Next come the field operators. There are a couple of technical complications. First, fields are like distributions. They need not be well-defined at points, but instead have to

be smeared against nice test functions. Secondly, even after smearing, fields are often unbounded operators. So there are domain questions. Gårding-Wightman handles these by postulating the existence of a common dense domain.

*2. Field Operators.* There exists a dense linear subspace  $D$  of  $F$  and for each  $f \in S(\mathbb{R}^4), j \in \{1, \dots, n\}$ , there exists an operator  $\phi_j(f)$  with domain containing  $D$  such that

$$\phi_j(f)D \subset D$$

$$\phi_j(f)^*D \subset D$$

$$f \rightarrow (\phi_j(f)v, w) \in S'(\mathbb{R}^4) \text{ for each } v, w \in D.$$

As with ordinary tempered distributions one writes

$$\phi_j(f) = \int dt d^3x \phi_j(t, x) f(t, x).$$

Time evolution is determined by a special self-adjoint operator, the Hamiltonian.

*3. Hamiltonian.* There exists a self-adjoint operator  $H$  such that

$$e^{isH}D \subset D \quad \forall s \in \mathbb{R}$$

$$e^{isH}\phi_j(t, x)e^{-isH} = \phi_j(t + s, x)$$

$$H \geq 0$$

$$0 \text{ is a simple eigenvalue of } H \text{ with eigenvector } \Omega \in D.$$

In practice the vector  $\Omega$  plays an important role. For example, in Quantum Electrodynamics,  $\Omega$  is the vacuum. That is, the state corresponding to a completely empty world. All physically interesting states are constructed by adding small numbers of particles to the vacuum.

In a relativistic QFT one wants covariance not just under time translations but also under the full Poincaré group. So there are axioms postulating the existence and basic properties of a unitary representation of (the covering group, inhomogeneous  $SL(2, \mathbb{C})$ , of) the Poincaré group. One also wants causality, i.e. measurements made at space-like separated points better not influence each other, so there is an axiom imposing this.

A constructive quantum field theorist is someone who tries, at least, to rigorously construct and determine properties of quantum field theories.

Just as in classical mechanics the Hamiltonian  $H = H_\lambda$  of a physically interesting system invariably contains a parameter (or several parameters)  $\lambda$  called the coupling constant. For a special value of  $\lambda$ , normally 0, the Hamiltonian is trivial in the sense that one can determine its properties and, in particular, the time evolution it generates, explicitly. So it is natural to view  $H_\lambda$  as a perturbation of  $H_0$ . However, for a QFT, this perturbation is extremely singular. I will discuss three symptoms of this singularity and how they are treated.

## 2. Change of Hilbert Space

The first appears even in the models that are easiest to deal with. The construction of the, so called, weakly coupled  $\lambda\phi_2^4$  model can now be reasonably presented in an advanced undergraduate analysis course. This model is completely characterized by a measure  $d\mu_\lambda$  on  $S'(R^2)$ . Even though not true, you might want to interpret  $S'(R^2)$  as the space of possible values of the field  $\phi(t, x)$  and the measure as the probability density for the field configurations when the coupling constant takes the value  $\lambda$ .

**Theorem [Fr].** *If  $0 \leq \lambda \neq \lambda'$  then  $d\mu_\lambda$  and  $d\mu_{\lambda'}$  are mutually singular.*

This theorem is not as much of an impediment as you might think. It suffices to deal with the measure weakly. In terms of the field operators, you should attempt to construct inner products  $(\Pi\phi_j(f_j)\Omega, \Omega)$  rather than to construct the field operators directly.

## 3. Renormalization

The second symptom is much more serious. It is the problem of renormalization [BS, H, FHRW]. I will illustrate it using the crystal model mentioned earlier [BG, FT]. For simplicity I will even discard the lattice, leaving only the “electrons”. Consider a system of  $N$  particles each of mass  $m$  living in a periodic box of side  $L$  and interacting with each other through an even two body potential  $\lambda V \in S(R^3)$ . The Hamiltonian for this system is

$$H_{L,N,\lambda} = \frac{1}{2m} \sum_{i=1}^N (-A_{x_i}) + \lambda \sum_{i < j} V(\mathbf{x}_i - \mathbf{x}_j)$$

and acts on the Hilbert space

$$F_N = \left\{ v \in L^2((R^3/LZ^3)^N) \mid v(\mathbf{x}_{\sigma(1)}, \dots, \mathbf{x}_{\sigma(N)}) = (-1)^{\text{sgn}(\sigma)} v(\mathbf{x}_1, \dots, \mathbf{x}_N) \text{ for all } \sigma \in S_N \right\}.$$

The antisymmetry of the vectors  $v \in F_N$  reflects the hypothesized Fermionic character of the particles. We would now like to take the thermodynamic limit  $L, N \rightarrow \infty$  with the density  $\varrho = N/L^3$  fixed. In order to avoid having to introduce the electron field I will concentrate on the ground state energy density

$$p(\lambda) = \lim_{\substack{L,N \rightarrow \infty \\ N/L^3 = \varrho}} \frac{\inf \text{spec } H_{L,N,\lambda}}{L^3}.$$

Of course it is not easy to determine  $\inf \text{spec } H_{L,N,\lambda}$  except when  $\lambda = 0$ . So we do some perturbation theory. Assume  $p(\lambda)$  exists (even this is not obvious), is  $C^\infty$  and that we can move derivatives inside the limit sign. Then it is not so difficult to develop an explicit formula for  $\frac{d^n p}{d\lambda^n}(0)$ . This formula is conventionally stated in terms of Feynman diagrams

$$\frac{d^n p}{d \lambda^n}(0) = \sum_{\substack{\text{connected vacuum} \\ \text{graphs of order } n}} \text{Val}(G).$$

The sum is over roughly  $(n!)^2$  Feynman diagrams each of which is a mnemonic device for a specific integral, called  $\text{Val}(G)$ .

For example one diagram with  $n = 2$  is given in Fig. 1.

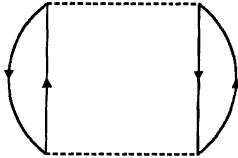


Fig. 1

It has the value

$$\begin{aligned} & \int \frac{d^4 p}{(2\pi)^4} \frac{d^4 q}{(2\pi)^4} \frac{d^4 r}{(2\pi)^4} \widehat{V}(\mathbf{q})^2 \frac{1}{ip_0 - e(\mathbf{p})} \frac{1}{i(p_0 - q_0) - e(\mathbf{p} - \mathbf{q})} \frac{1}{i(r_0 - q_0) - e(\mathbf{r} - \mathbf{q})} \\ & \quad \times \frac{1}{ir_0 - e(\mathbf{r})} \quad \text{where } p = (p_0, \mathbf{p}), \ q = (q_0, \mathbf{q}), \ r = (r_0, \mathbf{r}), \\ & \quad e(\mathbf{p}) = \frac{1}{2m} \mathbf{p}^2 - \mu_0 \\ & \quad \mu_0 = \frac{\pi}{4m} \left(\frac{3}{2}\sqrt{\pi}\right)^{3/2}. \end{aligned}$$

It is not clear whether the integral giving the value of this (or for that matter any other) graph converges. There are two potential obstructions to convergence. Firstly the integrand decays rather slowly at infinity. This turns out to be harmless. All graphs turn out to be convergent (though not necessarily absolutely convergent) at infinity. Secondly the integrand has singularities. The factor  $ip_0 - e(\mathbf{p})$  is zero when  $p_0 = 0$  and  $|\mathbf{p}| = \sqrt{2m\mu_0}$ . More precisely

$$\frac{\text{const}'}{\sqrt{p_0^2 + R^2}} \leq \left| \frac{1}{ip_0 - e(\mathbf{p})} \right| \leq \frac{\text{const}}{\sqrt{p_0^2 + R^2}}, \quad R = |\mathbf{p}| - \sqrt{2m\mu_0}$$

on  $|\mathbf{p}| \leq \text{const}$ . Hence  $[ip_0 - e(\mathbf{p})]^{-1}$  is locally  $L^1$  but not locally  $L^n$  for any  $n \geq 2$ .

While the graph in the example converges there are tons of them, in fact infinitely many, that diverge. It turns out that, in this model, every single divergent graph has

$$\text{Val}(G) = \int \frac{d^4 k}{(2\pi)^4} \frac{f(k)}{[ik_0 - e(\mathbf{k})]^n} \quad n \geq 2$$

with  $n \geq 2$ . I have lumped all the other integrals into  $f(k)$ . These integrals may also diverge if they have this same form.

I claim that all of these divergences arise, not because of any pathology in the model, but merely because we are trying to do the perturbation expansion in a stupid way. To see what I mean let's pretend that the exact

$$p(\lambda) = \int \frac{d^4 k}{(2\pi)^4} \frac{f(k, \lambda)}{ik_0 - \frac{1}{2m}\mathbf{k}^2 + \mu_0 + \delta\mu(\mu_0, \lambda)}$$

for some reasonable functions  $f(k, \lambda)$  and  $\delta\mu(\mu_0, \lambda)$  with  $\delta\mu(\mu_0, 0) = 0$ . The integral converges. Yet when we apply  $(d^n/d\lambda^n)|_{\lambda=0}$  we get large powers of  $[ik_0 - e(\mathbf{k})]^{-1}$  and hence divergence. It is not wise to try to expand a function,  $[ik_0 - e(\mathbf{k}) + \delta\mu]^{-1}$ , with a singularity somewhere, in powers of a function,  $[ik_0 - e(\mathbf{k})]^{-1}$  with a singularity somewhere else.

Once we know that this is the source of the divergences, the treatment is simple. Just parametrize our models, not by  $\mu_0$  and  $\lambda$ , but rather by  $\mu = \mu_0 + \delta\mu(\mu_0, \lambda)$  and  $\lambda$ . Then the position of the singularity in

$$\frac{1}{ip_0 - e(\mathbf{p}) + \delta\mu(\mu_0, \lambda)} = \frac{1}{ip_0 - \frac{1}{2m}\mathbf{p}^2 + \mu}$$

becomes independent of  $\lambda$  and the source of divergence mentioned above disappears. Of course to implement this you still have to find the right  $\delta\mu(\mu_0, \lambda)$ . But there is a renormalization algorithm which tells you how to choose  $\delta\mu(\mu_0, \lambda)$ , inductively in powers of  $\lambda$ .

**Theorem [FT].** *There exists a formal power series (with finite coefficients),  $\mu_0(\lambda, \mu) \sim \mu + \sum_{i=1}^{\infty} m_i(\mu) \lambda^i$ , such that every coefficient in the composite f.p.s.*

$$p(\lambda, \mu_0(\lambda, \mu), L) \sim \sum_{n=0}^{\infty} p_n^{\text{ren}}(L) \lambda^n$$

converges as  $L \rightarrow \infty$ . Furthermore

$$|\lim_{L \rightarrow \infty} p_n^{\text{ren}}(L)| \leq \text{const}^n n!.$$

This phenomenon – the divergence of coefficients in a perturbation expansion unless parameters are carefully adjusted so as to hold some physical quantity fixed – occurs commonly. For example, in classical mechanics, when you perturb a periodic orbit in a system of one degree of freedom you must be careful to keep the period fixed. That is, it is not generally possible to find a symplectic map converting  $\frac{1}{2}(p^2 + q^2) + \lambda v(p, q)$  to  $\frac{1}{2}(P^2 + Q^2)$ . One must use  $f_\lambda(\frac{1}{2}(P^2 + Q^2))$  with a very carefully chosen  $f_\lambda$  as the target.

## 4. Symmetry Breaking

The final symptom of the singularity of perturbations in QFT that I will discuss is the phenomenon of symmetry breaking. It is possible for all the input data used to specify a model to have a given symmetry without the resulting model having the symmetry. At first this may seem like a shocking, even inflammatory, statement. It is possible to rephrase it in an innocuous form but I feel that the above one gives a more realistic reflection of the impact of symmetry breaking.

Once again my illustration uses the many-electron model [FW, S]. Recall that the space of  $N$ -electron states in our periodic box is

$$F_N = \{v\varepsilon L^2((R^3/LZ^3)^N) | v(\mathbf{x}_{\sigma(1)}, \dots, \mathbf{x}_{\sigma(N)}) = (-1)^\sigma v(\mathbf{x}_1, \dots, \mathbf{x}_N) \quad \forall \sigma \in S_N\}.$$

To avoid having to change Hilbert spaces every time we change  $N$  we consider

$$F = \bigoplus_N F_N$$

with the modified “Hamiltonian”

$$\tilde{H}_L = H_{L,\lambda} - \mu N - c(L)1$$

( $c(L)$  is chosen to make  $\inf \text{spec } \tilde{H}_L = 0$ ) defined by

$$\tilde{H}_L|F_N = \frac{1}{2m} \sum_{i=1}^N -\Delta_{\mathbf{x}_i} + \sum_{i < j} \lambda V(\mathbf{x}_i - \mathbf{x}_j) + \text{interaction with ions} - \mu_0 N - c(L)1.$$

As  $L \rightarrow \infty$ , states with the wrong density end up with “ $\tilde{H} = \infty$ ” and get pushed out of the domain of  $\tilde{H}$ .

Now

$$V(\alpha) = e^{i\alpha(N-d(L))}$$

( $d(L)$  is a normalization constant) is a unitary representation of  $R$ , the covering group of  $U(1)$ , which commutes with  $\tilde{H}_L$  for every  $L$  since  $\tilde{H}_L : F_N \rightarrow F_N$ . Suppose that  $\tilde{H} = \lim_{L \rightarrow \infty} \tilde{H}_L$  were some ordinary (e.g. strong) limit. If  $\tilde{H}$  has a simple eigenvalue at the bottom of its spectrum, then the corresponding eigenvector  $\Omega$  must also be an eigenvector of  $V(\alpha)$  and in fact

$$V(\alpha)\Omega = \Omega$$

if the normalization constant  $d(L)$  is chosen appropriately.

The value  $(O\Omega, \Omega)$  of the observable (= operator)  $O$  measured in the state  $\Omega$  obeys  $(O\Omega, \Omega) = (V(\alpha)^* O V(\alpha)\Omega, \Omega)$ . There are whole orbits of observables that give the same answer when measured in the state  $\Omega$ . The symmetry  $U(1)$  is not broken.

On the other hand if there is an eigenspace of dimension greater than one at the bottom of the spectrum of  $\tilde{H}$ , then each individual eigenvector  $\Omega$  need not be an eigenvector of  $V(\alpha)$ , though of course the eigenspace is still invariant under  $V(\alpha)$ . The value of  $O$  measured in the state  $\Omega$  need no longer be constant as  $O$  moves over the orbit.

In a QFT the limit  $\tilde{H} = \lim_{L \rightarrow \infty} \tilde{H}_L$  is not an ordinary (e.g. strong) limit. Instead of getting one model with a degenerate ground state one gets a whole

family [I, R] of independent models  $F^{(\beta)}$  each with a unique ground state. Each independent model no longer carries a representation of  $U(1)$  because  $V(\alpha)$  tries to map  $\Omega_\beta \in F^{(\beta)}$  to  $\Omega_{\alpha+\beta} \in F^{(\alpha+\beta)} \neq F^{(\beta)}$ .

This is symmetry breaking. It has important physical consequences. In the example we have been discussing superconductivity arises. It also has a big impact on any attempted construction. In practice one must always attempt to construct the interacting broken symmetry model as a perturbation of a free model which also has the same broken symmetry.

## References

- [BS] Bogoliubov, N.N., Shirkov, D.V.: Introduction to the theory of quantized fields. Wiley, New York 1980
- [BG] Benfatto, G., Gallavotti, G.: Perturbation theory of the Fermi surface in a quantum liquid. A general quasi particle formalism and one dimensional systems. *J. Stat. Phys.* (1990) (to appear)
- [FHRW] Feldman, J., Hurd, T., Rosen, L., Wright, J.: QED: A proof of renormalizability. (Lecture Notes in Physics, vol. 312.) Springer, Berlin Heidelberg New York 1988
- [FT] Feldman, J., Trubowitz, E.: Perturbation theory for many fermion systems. *Helv. Phys. Acta* **63** (1990) 157–260
- [FW] Fetter, A.L., Walecka J.D.: Quantum theory of many-particle systems. McGraw Hill, New York 1971
- [Fr] Froehlich, J.: Verification of axioms for Euclidean and relativistic fields and Haag's theorem in a class of  $P(\phi)_2$  models. *Ann. L'Inst. H. Poincaré* **21** (1974) 271–317
- [H] Hepp, K.: Théorie de la renormalisation. (Lecture Notes in Physics, vol. 2.) Springer, Berlin Heidelberg New York 1969
- [I] Israel, R.B.: Convexity in the theory of lattice gases. Princeton University Press, Princeton 1979
- [R] Ruelle, D.: Statistical mechanics, rigorous results. Benjamin, New York Amsterdam 1969
- [S] Schrieffer, J.R.: Theory of superconductivity. (Frontiers in Physics, vol. 20.) Addison-Wesley, Redwood City 1964
- [SW] Streater, R., Wightman, A.: PCT, spin and statistics and all that. Benjamin, New York Amsterdam 1964



# Solvable Lattice Models and Quantum Groups

Michio Jimbo

Department of Mathematics, Faculty of Science, Kyoto University, Kyoto 606, Japan

Our main concern in this survey is the following question. What is the rôle of quantum groups in solvable lattice models; more specifically, what is their rôle in building models and in solving models respectively?

Building models amounts to constructing solutions to the Yang-Baxter equation (YBE). Through the development of the last 4–5 years, the situation has been clarified relatively well. In the context of quantum groups, the theory of YBE is the theory of intertwiners of their representations. In the first half of this article we shall elucidate this statement on three constructions.

In contrast, the rôle of quantum groups in solving models is very obscure. By ‘solving’ the model thus constructed we mean to calculate physically important quantities in closed form. The one point functions are one of the few quantities which can be evaluated exactly, thanks to Baxter’s corner transfer matrix method. In many cases they are expressible in terms of certain modular forms arising from representations of affine Lie algebras. In the latter half of this paper we shall explain these results in the light of Kashiwara’s theory of crystal base for quantum groups.

## 1. Yang-Baxter Equation

Let  $V$  be a finite dimensional vector space over  $\mathbf{C}$ , and let  $R(\xi, \eta)$  be a function of some variables  $\xi, \eta$  with values in  $\text{End}(V \otimes V)$ . Then YBE is the following functional equation for  $R(\xi, \eta)$ , written in  $\text{End}(V \otimes V \otimes V)$ :

$$R_2(\xi, \eta)R_1(\xi, \zeta)R_2(\eta, \zeta) = R_1(\eta, \zeta)R_2(\xi, \zeta)R_1(\xi, \eta) \quad (1)$$

where

$$R_1(\xi, \eta) = R(\xi, \eta) \otimes 1, \quad R_2(\xi, \eta) = 1 \otimes R(\xi, \eta).$$

Following the common terminology we call a solution of YBE an  $R$  matrix, and  $\xi, \eta, \dots$  spectral parameters.

Throughout this article, by quantum groups we mean the quantized enveloping algebras  $U_q(\mathfrak{g})$ . For us the most important case is when  $\mathfrak{g}$  is an *affine* (rather than finite dimensional) Lie algebra. Let us fix notations.  $(a_{ij})_{1 \leq i, j \leq l}$  signifies the Cartan matrix of  $\mathfrak{g}$ ,  $d_i \in \mathbf{Q}^\times$  are such that  $d_i a_{ij} = d_j a_{ji}$ , and  $q$  is a complex

number satisfying  $q_i^2 \neq 0$  with  $q_i = q^{d_i}$ . Then  $U_q(\mathfrak{g})$  is the  $\mathbf{C}$ -algebra with the generators  $e_i, f_i, t_i^{\pm 1}$  ( $1 \leq i \leq l$ ) subject to the following defining relations.

$$\begin{aligned} t_i t_j &= t_j t_i, \quad t_i t_i^{-1} = 1 = t_i^{-1} t_i, \\ t_i e_j t_i^{-1} &= q_i^{a_{ij}} e_j, \quad t_i f_j t_i^{-1} = q_i^{-a_{ij}} f_j, \quad [e_i, f_j] = \delta_{ij} \frac{t_i - t_i^{-1}}{q_i - q_i^{-1}}, \\ \sum_{l=0}^{1-a_{ij}} (-1)^l e_i^{(1-a_{ij}-l)} e_j e_i^{(l)} &= 0, \quad \sum_{l=0}^{1-a_{ij}} (-1)^l f_i^{(1-a_{ij}-l)} f_j f_i^{(l)} = 0, \quad (i \neq j) \end{aligned}$$

where  $e_i^{(l)} = e_i^l / [l]_{q_i}!$ ,  $f_i^{(l)} = f_i^l / [l]_{q_i}!$ ,  $[l]_t! = \prod_{j=1}^l (t^j - t^{-j}) / (t - t^{-1})$ . The comultiplication is given by

$$\Delta(e_i) = e_i \otimes 1 + t_i \otimes e_i, \quad \Delta(f_i) = f_i \otimes t_i^{-1} + 1 \otimes f_i, \quad \Delta(t_i) = t_i \otimes t_i.$$

## 1.1 Trigonometric $R$ Matrices

The first construction associates an  $R$  matrix with each pair  $(U_q(\mathfrak{g}), \pi)$ , consisting of a quantum group and its finite dimensional irreducible representation. Here  $q$  is assumed to be ‘generic’, i. e.  $q^n \neq 1$  for all positive integers  $n$ .

Let us take the simplest example of  $U_q(\widehat{\mathfrak{sl}}(2, \mathbf{C}))$  corresponding to  $(a_{ij}) = \begin{pmatrix} 2 & -2 \\ -2 & 2 \end{pmatrix}$ , and its two dimensional representation  $(\pi_\xi, V = \mathbf{C}^2)$  depending on  $\xi \in \mathbf{C}^\times$  given as follows.

$$\begin{aligned} \pi_\xi(e_0) &= \begin{pmatrix} 0 & 0 \\ \xi & 0 \end{pmatrix}, \quad \pi_\xi(f_0) = \begin{pmatrix} 0 & \xi^{-1} \\ 0 & 0 \end{pmatrix}, \quad \pi_\xi(t_0) = \begin{pmatrix} q^{-1} & 0 \\ 0 & q \end{pmatrix}, \\ \pi_\xi(e_1) &= \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, \quad \pi_\xi(f_1) = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \quad \pi_\xi(t_1) = \begin{pmatrix} q & 0 \\ 0 & q^{-1} \end{pmatrix}. \end{aligned}$$

An intertwiner  $R(\xi, \eta)$  between  $\pi_\xi \otimes \pi_\eta$  and  $\pi_\eta \otimes \pi_\xi$  is a linear isomorphism  $V \otimes V \rightarrow V \otimes V$  such that

$$R(\xi, \eta) (\pi_\xi \otimes \pi_\eta) \circ \Delta(X) = (\pi_\eta \otimes \pi_\xi) \circ \Delta(X) R(\xi, \eta) \quad \forall X \in U_q(\widehat{\mathfrak{sl}}(2, \mathbf{C})).$$

Solving these linear equations for  $R$  one finds that up to a scalar

$$R(\xi, \eta) = \begin{pmatrix} a & 0 & 0 & 0 \\ 0 & c & b & 0 \\ 0 & b & \bar{c} & 0 \\ 0 & 0 & 0 & a \end{pmatrix},$$

with  $a = \xi q - \eta q^{-1}$ ,  $b = \xi - \eta$ ,  $c = \xi(q - q^{-1})$  and  $\bar{c} = \eta(q - q^{-1})$ . This is the well known  $R$  matrix of the 6 vertex model [1].

The left and right hand sides of YBE (1) are both intertwiners of the threefold tensor products  $\pi_\xi \otimes \pi_\eta \otimes \pi_\zeta \rightarrow \pi_\zeta \otimes \pi_\eta \otimes \pi_\xi$ . It can be checked that, if the spectral

parameters are generic, any intertwiner of  $\pi_\xi \otimes \pi_\eta \otimes \pi_\zeta$  into itself must be a scalar. From this one deduces that  $R(\xi, \eta)$  solves YBE.

The explicit form of the  $R$  matrices are available for the vector representation of  $U_q(\mathfrak{g})$  of classical types [2, 3]; further results are given in [4–8]. In the general case the existence of the intertwiners is guaranteed by Drinfeld's universal  $R$  matrix [9]. Irreducible finite dimensional representations of Yangians, which are certain degeneration of  $U_q(\mathfrak{g})$ , are classified also by Drinfeld [10] (the results for  $U_q(\mathfrak{g})$  are almost the same; cf. [11]). In order to describe the corresponding  $R$  matrices in full generality, further work remains to be done.

## 1.2 Chiral Potts-Type Models

The second construction is related to  $R$  matrices whose spectral parameters live on certain algebraic curves. Here we take  $q$  to be a primitive  $N$ -th root of unity:  $q^N = 1$  (we assume  $N$  is odd). For our purposes we enlarge the algebra  $U_q(\mathfrak{g})$  slightly by adding central elements  $z_i^{\pm 1}$ , which modify the comultiplication as

$$\begin{aligned}\Delta(e_i) &= e_i \otimes 1 + z_i t_i \otimes e_i, & \Delta(f_i) &= f_i \otimes t_i^{-1} + z_i^{-1} \otimes e_i, \\ \Delta(t_i) &= t_i \otimes t_i, & \Delta(z_i) &= z_i \otimes z_i.\end{aligned}$$

The resulting algebra  $\tilde{U}_q(\mathfrak{g})$  is known as the quantum double of a Borel subalgebra of  $U_q(\mathfrak{g})$ . A characteristic feature of  $q$  being a root of 1 is that the elements  $e_i^N, f_i^N, t_i^N$  lie in the center  $\mathcal{Z}$  of  $\tilde{U}_q(\mathfrak{g})$ .

Again we take the example of  $\mathfrak{g} = \widehat{\mathfrak{sl}}(2, \mathbf{C})$ . We consider an  $N$  dimensional irreducible representation of  $\tilde{U}_q(\mathfrak{g})$  containing 5 parameters  $\xi = (a_0, a_1, x_0, x_1, c) \in (\mathbf{C}^\times)^5$ .

$$\begin{aligned}\pi_\xi(e_0) &= x_0 F, & \pi_\xi(f_0) &= x_0^{-1} E, & \pi_\xi(t_0) &= T^{-1}, \\ \pi_\xi(e_1) &= x_1 E, & \pi_\xi(f_1) &= x_1^{-1} F, & \pi_\xi(t_1) &= T, \\ \pi_\xi(z_0) &= c_0, & \pi_\xi(z_1) &= c_0^{-1}.\end{aligned}$$

Here

$$E = \frac{a_1^2 Z - 1}{q - q^{-1}} X, \quad F = (a_0 a_1 X)^{-1} \frac{a_0^2 Z^{-1} - 1}{q - q^{-1}}, \quad T = q \frac{a_1}{a_0} Z,$$

and

$$X = \begin{pmatrix} 1 & & & & \\ & q^2 & & & \\ & & \ddots & & \\ & & & q^{2N-2} & \end{pmatrix}, \quad Z = \begin{pmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ 1 & & & & 0 \end{pmatrix}.$$

Suppose that  $\pi_\xi \otimes \pi_\eta$  is equivalent to  $\pi_\eta \otimes \pi_\xi$ . Since  $\Delta(e_i^N) = e_i^N \otimes 1 + (z_i t_i)^N \otimes e_i^N \in \mathcal{Z} \otimes \mathcal{Z}$ , one must have

$$\frac{\pi_\xi(e_i)^N}{1 - \pi_\xi(z_i t_i)^N} = \frac{\pi_\eta(e_i)^N}{1 - \pi_\eta(z_i t_i)^N} (\equiv \Gamma_i)$$

and likewise for the  $f_i^N$ . Hence  $\xi, \eta$  are forced to lie on a common algebraic variety  $\mathcal{S}_\Gamma$  with  $\Gamma = (\Gamma_i)$  being the moduli parameters. It turns out that  $\mathcal{S}_\Gamma$  is essentially a product  $\mathcal{C}_k \times \mathcal{C}_k$  of curves

$$\mathcal{C}_k = \left\{ (x, y, \mu) \mid x^N + y^N = k(1 + x^N y^N), \quad \mu^N = \frac{\sqrt{1 - k^2}}{1 - kx^N} \right\}.$$

Conversely, if  $\xi = (r, r')$ ,  $\eta = (\tilde{r}, \tilde{r}')$  ( $r = (x, y, \mu) \in \mathcal{C}_k$  etc.) lie on  $\mathcal{S}_\Gamma$ , then there exists an intertwiner and is given by the following formulas.

$$\begin{aligned} R(\xi, \eta) &= R'(\tilde{r}, r') (R''(r', \tilde{r}') \otimes R''(r, \tilde{r})) R'(r, \tilde{r}'), \\ R'(r, \tilde{r}) &= \sum_{a=0}^{N-1} \widehat{W}_{r\tilde{r}}(a) (X^{-1} \otimes X)^a, \quad R''(r, \tilde{r}) = \sum_{a=0}^{N-1} \overline{W}_{r\tilde{r}}(a) Z^a, \\ \widehat{W}_{r\tilde{r}}(a) &= \prod_{l=1}^a \frac{\tilde{\mu}y - \mu\tilde{y}q^{2l-2}}{\tilde{\mu}\tilde{x} - \mu x q^{2l}}, \quad \overline{W}_{r\tilde{r}}(a) = \prod_{l=1}^a \mu\tilde{\mu} \frac{xq^2 - \tilde{x}q^{2l}}{\tilde{y} - yq^{2l}}. \end{aligned}$$

The same reasoning as before shows that  $R(\xi, \eta)$  solves YBE. The  $\widehat{W}_{rr}$ ,  $\overline{W}_{rr}$  are the Boltzmann weights (the former being Fourier transformed) of the chiral Potts model found by Au-Yang, Baxter, McCoy, Perk and others [12–13]. The connection with the quantum groups at roots of 1 was first noticed by Bazhanov-Stroganov [14] in a different language. The formulation above follows [15].

The study of intertwiners at roots of 1 has begun only recently. The chiral Potts model has been generalized to the case corresponding to a class of representations of  $U_q(\widehat{\mathfrak{sl}}(n, \mathbf{C}))$  [16, 17]. In the general case the structure of irreducible representations at  $q^N = 1$  and the existence of intertwiners are the major open questions. We remark that, for simple finite dimensional  $\mathfrak{g}$ , irreducible representations have been classified by De Concini and Kac [18].

### 1.3 Face Models

The third construction deals with solutions of YBE (1) which are not of the form  $R(\xi, \eta) \in \text{End}(V \otimes V)$ .

Let us return to the example of Sect. 1.1.,  $q$  being generic again. Let  $V_a$  denote the  $a$  dimensional irreducible representation of the subalgebra  $U_q(\mathfrak{sl}(2, \mathbf{C})) \subset U_q(\widehat{\mathfrak{sl}}(2, \mathbf{C}))$ . We consider the decomposition with respect to  $U_q(\mathfrak{sl}(2, \mathbf{C}))$  of the tensor product representation

$$\begin{aligned} V_a \otimes V_2^{\otimes N} &= \bigoplus_c \Omega_{ac}^N \otimes V_c, \\ \Omega_{ac}^N &= \{v \in V_a \otimes V_2^{\otimes N} \mid e_1 v = 0, \quad t_1 v = q^{c-1} v\}. \end{aligned}$$

The  $R = R(\xi, \eta) \in \text{End}(V_2 \otimes V_2)$  commutes with the action of  $U_q(\mathfrak{sl}(2, \mathbf{C}))$ , hence

$$W_i = \text{id} \otimes \text{id} \otimes \cdots \otimes R^{\otimes i+1} \otimes \cdots \otimes \text{id} \in \text{End}(V_a \otimes V_2^{\otimes N})$$

are well-defined operators on  $\Omega_{ac}^N$  and satisfy YBE.

Repeated use of the Clebsch-Gordan decomposition gives rise to a basis of  $\Omega_{ac}^{\mathcal{N}}$  indexed by a sequence  $\mu = (\mu_0, \mu_1, \dots, \mu_N)$  ( $\mu_0 = a, \mu_N = c$ ) such that  $\mu_i = \mu_{i+1} \pm 1$  for each  $i$ . In this basis the matrix elements of  $W_i$  look as follows.

$$(W_i)_{\mu v} = \delta_{\mu_0, v_0} \cdots \delta_{\mu_{i-1}, v_{i-1}} W \begin{pmatrix} \mu_{i-1} & v_i \\ \mu_i & \mu_{i+1} \end{pmatrix} \delta_{\mu_{i+1}, v_{i+1}} \cdots,$$

$$W \begin{pmatrix} a & a \pm 1 \\ a \pm 1 & a \pm 2 \end{pmatrix} = \frac{[1+u]}{[1]}, \quad W \begin{pmatrix} a & a \pm 1 \\ a \pm 1 & a \end{pmatrix} = \frac{[a \mp u]}{[a]},$$

$$W \begin{pmatrix} a & a \mp 1 \\ a \pm 1 & a \end{pmatrix} = \frac{[u]}{[1]} \sqrt{\frac{[a+1][a-1]}{[a]^2}},$$

where  $q^u = \xi/\eta$  and  $[u] = (q^u - q^{-u})/(q - q^{-1})$ . Starting from the Verma module instead of  $V_a$ , we obtain the same formulas with  $a$  being any complex number.

In fact if one replaces the symbol  $[u]$  by the elliptic theta function

$$\theta_1 \left( \frac{\pi u}{L}, p \right) = 2p^{1/8} \sin \frac{\pi u}{L} \prod_{k=1}^{\infty} \left( 1 - 2p^k \cos \frac{2\pi u}{L} + p^{2k} \right) (1 - p^k),$$

these formulas still solve YBE ( $L \in \mathbf{C}^\times$  is an arbitrary parameter). Such elliptic solutions are known in correspondence to many of the representations of  $U_q(\mathfrak{g})$  [19–23]. Pasquier [24] pointed out that their trigonometric degenerations ( $p \rightarrow 0$ ) can be described as above in terms of quantum groups. It remains an open problem to explain the existence of elliptic solutions.

## 2. One-Point Functions

### 2.1 Formulation

Let us now come to our second topic about the calculation of physical quantities called one point functions. The problem is set up as follows. We consider a square lattice  $\mathcal{L}$  on the plane  $\mathbf{Z}^2$ . Each site  $i = (i_1, i_2)$  of the lattice has a random variable  $s_i$  taking values in some set  $\mathcal{S}$ . A configuration is an assignment  $s = (s_i)$  of the values for each  $i$ . One also imposes the condition that the  $s_i$  for the boundary sites  $i \in \partial\mathcal{L}$  are fixed to a given configuration  $\bar{s}$  (which is to be chosen a ‘ground state’; see below).

A model is specified by giving a Boltzmann weight  $W \begin{pmatrix} a & b \\ d & c \end{pmatrix}$  for each quadruple  $(a, b, c, d)$  of elements of  $\mathcal{S}$ . Given a configuration  $s$ , to each ‘face’ (= an elementary square) of the lattice we attach a number  $W \begin{pmatrix} s_i & s_j \\ s_l & s_k \end{pmatrix}$  where  $i, j, k, l$  are the NW, NE, SE, SW corner sites of the face. The probability of a given configuration  $s$  to occur is defined to be

$$p(s) = Z^{-1} \prod_{\text{faces}} W \begin{pmatrix} s_i & s_j \\ s_l & s_k \end{pmatrix}$$

where the product ranges over the faces of  $\mathcal{L}$ , and  $Z$  is a normalization constant to make the total probability = 1. Let us now fix  $\bar{s}$ , a site  $0 \in \mathcal{L}$  and  $a \in \mathcal{S}$ . By definition the one point function is the expected value of  $\delta_{s_0, a}$ :

$$P_a(\bar{s}) = \text{Prob}(s_0 = a) = \lim \sum_s \delta_{s_0, a} p(s).$$

Here the lim signifies the limit of lattice size tending to  $\infty$ .

Let us take again the  $R$  matrix of Sect. 1.1. Taking base vectors  $\{v_\lambda\}_{\lambda=\pm 1} (\{\pm 1\} = \text{the set of weights of } V_2)$ , we set  $R(\xi, \eta)v_\lambda \otimes v_\mu = \sum v_\kappa \otimes v_\nu R(\xi, \eta)_{\kappa\nu, \lambda\mu}$ . Let  $\mathcal{S} \simeq \mathbf{Z}$  be the set of level 1 integral weights of  $\widehat{\mathfrak{sl}}(2, \mathbf{C})$ . For  $a, b, c, d \in \mathcal{S}$  we set

$$\begin{aligned} W \begin{pmatrix} a & b \\ d & c \end{pmatrix} &= 0 \quad \text{unless } a - b, b - c, a - d, d - c \text{ all belong to } \{\pm 1\}, \\ &= R(\xi, \eta)_{\kappa\nu, \lambda\mu} \quad \text{if } a + \lambda = b, b + \mu = c, a + \kappa = d, d + \nu = c. \end{aligned}$$

Hereafter we consider the region  $0 < q < 1 < x = \xi/\eta$ .

The boundary configuration  $\bar{s}$  is so chosen as to maximize  $p(s)$ . In the present case,  $\bar{s}$  turns out to be in one-to-one correspondence with dominant integral weights of fixed level ( $= 1$  here). For example we take  $\bar{s}_i = 0$  or 1 according as  $i_1 + i_2$  is even or odd. This corresponds to the fundamental weight  $A_0$ .

Baxter's corner transfer matrix (CTM) is the only known tool to evaluate the one point function exactly. It is an infinite dimensional matrix whose trace gives the one point function up to a simple factor. In the sequel we shall focus on this trace. According to Baxter's argument based on YBE [1], it is completely determined by the trace at  $q = 0$ , where CTM becomes diagonal. We quote here the conclusion of the method. The CTM eigenvectors at  $q = 0$  are labeled by an object called *path*. In the present example, it is a sequence  $\mu = (\mu_0, \mu_1, \mu_2, \dots)$  ( $\mu_i \in \mathcal{S}$ ) such that

$$\mu_i = \mu_{i-1} \pm 1 \quad \text{for all } i, \quad \mu_i = \bar{\mu}_i \quad \text{for } i \gg 0,$$

where  $\bar{\mu}_i = (0, 1, 0, 1, 0, 1, \dots)$  corresponds to the ground state  $\bar{s}$ . The trace of CTM becomes the combinatorial sum (called the one dimensional configuration sum)

$$\begin{aligned} c_a^{A_0}(t) &= \sum_{\mu} \delta_{\mu_0, a} t^{\omega(\mu)}, \\ \omega(\mu) &= \sum_{j=1}^{\infty} j (H(\mu_{j-1}, \mu_j, \mu_{j+1}) - H(\bar{\mu}_{j-1}, \bar{\mu}_j, \bar{\mu}_{j+1})). \end{aligned} \tag{2}$$

Here  $H(\lambda, \mu, \nu) = 0$  if  $\lambda + 1 = \mu = \nu + 1, = 0$  otherwise, and the sum ranges over the paths whose 'initial point'  $\mu_0$  is fixed to  $a \in \mathcal{S}$ . This is a counting problem – counting the number of paths such that  $\mu_0$  and  $\omega(\mu)$  are fixed. It is not hard to show that

$$c_a^{A_0}(t) = \frac{t^{a^2/4}}{\prod_{j=1}^{\infty} (1 - t^j)}.$$

The quantity on the right hand side coincides with the string function of the irreducible highest weight module  $M(\Lambda_0)$  of the affine Lie algebra  $\widehat{\mathfrak{sl}}(2, \mathbf{C})$  in the sense of [25].

## 2.2 Crystal Base and One-Point Functions

In the same manner each  $R$  matrix of type  $(U_q(\mathfrak{g}), \pi)$  gives rise to a similar model. The theory of crystal base [26] provides a natural framework to deal with the sum (2) in general. According to [26], for any integrable highest weight representation  $M(\Lambda)$  of  $U_q(\mathfrak{g})$  with highest weight  $\Lambda$ , there exists a unique canonical base  $B(\Lambda)$  (crystal base) ‘at  $q = 0$ ’. Moreover they behave extremely simply with respect to taking tensor products. For details we refer the reader to the exposition of Kashiwara [27] in these proceedings.

Let us formulate the results in a generalized setting. We take

$$U_q(\mathfrak{g}) = U_q(\widehat{\mathfrak{sl}}(n, \mathbf{C})),$$

$V = l$ -th symmetric tensors  $S^l(\mathbf{C}^n)$  of  $n$  dimensional representation of  $U_q(\mathfrak{g})$ ,

$\mathcal{S}$  = the set of integral weights of level  $l$ .

Let  $\Lambda_i$  denote the fundamental weights of  $\widehat{\mathfrak{sl}}(n, \mathbf{C})$  and let  $\varepsilon_i = \Lambda_{i+1} - \Lambda_i$ . The set of weights of  $V$  is identified with  $\mathcal{A}_l = \{\varepsilon_{i_1} + \cdots + \varepsilon_{i_l} \mid 0 \leq i_1, \dots, i_l \leq n-1\}$ . Let  $P_l^+$  denote the set of dominant integral weights of level  $l$ . For  $\Lambda \in P_l^+$ , a  $\Lambda$ -path is defined to be a sequence  $\mu = (\mu_0, \mu_1, \mu_2, \dots)$  ( $\mu_i \in \mathcal{S}$ ) such that

$$\mu_i - \mu_{i-1} \in \mathcal{A}_l \text{ for all } i, \quad \mu_i = \sigma^i(\Lambda) \text{ for } i \gg 0,$$

where  $\sigma$  signifies the linear automorphism such that  $\sigma(\Lambda_i) = \Lambda_{i+1}$ . Finally if  $\mu - \lambda = \varepsilon_{i_1} + \cdots + \varepsilon_{i_l}$  and  $\nu - \mu = \varepsilon_{j_1} + \cdots + \varepsilon_{j_l}$ , then

$$H(\lambda, \mu, \nu) = \max_{\tau \in \mathcal{C}_l} \sum_{k=1}^l \theta(i_{\tau(k)} - j_k) \quad (3)$$

with  $\theta(x) = 0$  if  $x < 0$ ,  $= 1$  if  $x \geq 0$ . With these notations we have

**Theorem 1** [28, 29]. *There exists a bijective correspondence between the set of  $\Lambda$ -paths and the crystal base  $B(\Lambda)$  of  $M(\Lambda)$ , such that each path  $\mu$  has weight  $\mu_0 - \omega(\mu)\delta$  in the latter. Here  $\delta$  signifies the null root.*

**Corollary 2.** *The one dimensional sum (2) is the string function*

$$c_a^\Lambda(t) = \sum_{n \in \mathbf{Z}} \dim M(\Lambda)_{a-n\delta} t^n.$$

This result has been found first by a purely combinatorial method in [30]. There is also a relative version of these results. Here one takes

$$\mathcal{S} = P_{l'+l}^+,$$

$$W \begin{pmatrix} a & b \\ d & c \end{pmatrix} = \text{elliptic solutions in Sect.1.3. corresponding to } S^l(\mathbf{C}^n).$$

In this case  $\bar{s}$  corresponds to a pair  $(\Lambda', \Lambda) \in P_{l'}^+ \times P_l^+$ . A sequence  $\mu = (\mu_0, \mu_1, \mu_2, \dots)$  is called a  $(\Lambda', \Lambda)$ -path if

- a) for any  $i$  there exists a pair  $(\tilde{\Lambda}', \tilde{\Lambda}) \in P_{l'}^+ \times P_l^+$  such that  $\mu_i = \tilde{\Lambda}' + \tilde{\Lambda}$  and  $\mu_{i+1} = \tilde{\Lambda}' + \sigma(\tilde{\Lambda})$ ,
- b)  $\mu_i = \Lambda' + \sigma^i(\Lambda)$  for  $i \gg 0$ .

**Theorem 3** [29]. *A vector  $u \otimes v \in B(\Lambda') \otimes B(\Lambda)$  is a highest weight vector if and only if  $u$  is highest and  $v$  corresponds to a  $(\Lambda', \Lambda)$ -path.*

Let

$$M(\Lambda') \otimes M(\Lambda) = \bigoplus_{a \in P_{l'+l}^+} \Omega_{\Lambda'\Lambda a} \otimes M(a)$$

be the decomposition of  $U_q(\widehat{\mathfrak{sl}}(n, \mathbf{C}))$  modules where  $\Omega_{\Lambda'\Lambda a}$  denotes the space of highest weight vectors of weight  $a$ . Let  $c_a^{\Lambda', \Lambda}$  denote the one dimensional configuration sum (3) where  $\mu$  ranges over  $(\Lambda', \Lambda)$ -paths such that  $\mu_0 = a$ , and  $H(\lambda, \mu, v)$  is the same as (3). Then we have

**Corollary 4.**

$$c_a^{\Lambda', \Lambda}(t) = \sum_n \dim(\Omega_{\Lambda'\Lambda a})_{a-n\delta} t^n.$$

The left hand side is called the branching function studied in [25, 31]. Similar results for other classical type algebras (for level 1) can be found in [6, 7, 32].

### 3. Summary

We have seen two aspects of the interplay between solvable lattice models and quantum groups. Our conclusions are summarized as follows:

- 1) The theory of intertwiners lead to solutions of YBE, and hence to construction of solvable lattice models.
- 2) The one point functions are reduced to one dimensional configuration sums by CTM method. The theory of crystal base enables one to identify the latter with modular forms arising in the representation theory of affine Lie algebras.

Although crystal base provides a powerful combinatorial tool, there is no direct explanation why these modular forms should arise at all. We feel the intrinsic understanding is still missing.

## References

1. Baxter, R. J.: Exactly Solved Models in Statistical Mechanics. Academic, London 1982
2. Bazhanov, V.V.: Integrable quantum systems and classical Lie algebras. Comm. Math. Phys. **113** (1987) 471–503
3. Jimbo, M.: Quantum  $R$  matrix for the generalized Toda system. Comm. Math. Phys. **102** (1986) 537–547
4. Reshetikhin, N. Yu.: The spectrum of the transfer matrices connected with Kac-Moody algebras. Lett. Math. Phys. **14** (1987) 235–246
5. Reshetikhin, N. Yu.: Zapiski Nauchnyi Seminarov LOMI **169** (1988) 122
6. Okado, M.: Quantum  $R$  matrices related to the spin representations of  $B_n$  and  $D_n$ . Comm. Math. Phys. **134** (1990) 467–486
7. Kuniba, A.: Quantum  $R$  matrix for  $G_2$  and a solvable 175 vertex models. J. Phys. A **23** (1990) 1349–1362
8. Chung, H. J., Koh, I. G.: Solutions to the quantum Yang-Baxter equation for the exceptional Lie Algebras with a spectral parameter. Preprint 1990
9. Drinfel'd, V. G.: Quantum groups. Proceedings of the International Congress of Mathematicians, Berkeley, California, USA 1986, pp. 798–820
10. Drinfel'd, V. G.: A new realization of Yangians and quantized affine algebras. Sov. Math. Dokl. **36** (1988) 212–216
11. Chari, V., Pressley, A.: Quantum affine algebras. Preprint 1990
12. Au-Yang, H., McCoy, B. M., Perk, J. H. H., Tang, S.: Solvable models in statistical mechanics and Riemann surfaces of genus greater than one. In: Algebraic Analysis, vol. 1, eds. M. Kashiwara and T. Kawai. Academic, 1989, pp. 29–39; and references therein
13. Baxter, R. J., Perk, J. H. H., Au-Yang, H.: New solutions of the star-triangle relations for the chiral Potts model. Phys. Lett. A **128** (1988) 138–142
14. Bazhanov, V. V., Stroganov, Yu. G.: Chiral Potts model as a descendant of the six vertex models. J. Stat. Phys. **51** (1990) 799–817
15. Date, E., Jimbo, M., Miki, K., Miwa, M.: New  $R$  matrices associated with cyclic representations of  $U_q(A_2^{(2)})$ . Preprint RIMS **706** (1990)
16. Bazhanov, V. V., Kashaev, R. M., Mangazeev, V. V., Stroganov, Yu. G.:  $(\mathbb{Z}_N \times)^{n-1}$  generalization of the chiral Potts model. Preprint 1990
17. Date, E., Jimbo, M., Miki, K., Miwa, M.: Generalized chiral Potts models and minimal cyclic representations of  $U_q(\widehat{\mathfrak{gl}}(n, \mathbb{C}))$ . Preprint RIMS **715** (1990), Comm. Math. Phys. (to appear)
18. De Concini, C., Kac, V. G.: Representations of quantum groups at roots of 1. Preprint 1990
19. Baxter, R. J.: Eight-vertex model in lattice statistics and one-dimensional anisotropic Heisenberg chain II. Equivalence to a generalized ice-type model. Ann. Phys. **76** (1973) 25–47
20. Andrews, G. E., Baxter, R. J. and Forrester, P. J.: Eight-vertex SOS model and generalized Rogers-Ramanujan-type identities. J. Stat. Phys. **35** (1984) 193–266
21. Date, E., Jimbo, M., Kuniba, A., Miwa, T., Okado, M.: Exactly solvable SOS models II : Proof of the star-triangle relation and combinatorial identities. Adv. Studies Pure Math. **16** (1988) 17–122
22. Jimbo, M., Miwa, T., Okado, M.: Solvable lattice models related to the vector representation of classical simple Lie algebras. Comm. Math. Phys. **116** (1988) 507–525
23. Jimbo, M., Kuniba, A., Miwa, T., Okado, M.: The  $A_n^{(1)}$  face models. Comm. Math. Phys. **119** (1988) 543–565
24. Pasquier, V.: Etiology of IRF models. Comm. Math. Phys. **118** (1988) 335–364

25. Kac, V. G., Peterson, D. H.: Infinite-dimensional Lie algebras, theta functions and modular forms. *Adv. Math.* **53** (1984) 125–264
26. Kashiwara, M.: Bases cristallines. *C. R. Acad. Sci. Paris* **311** (1990) 277–280
27. Kashiwara, M.: Crystallizing the  $q$  analogue of universal enveloping algebras. In these Proceedings, pp. 791–797
28. Misra, K. C., Miwa, T.: Crystal bases for the basic representations of  $U_q(\widehat{\mathfrak{sl}}(n, \mathbf{C}))$ . *Comm. Math. Phys.* **134** (1990) 79–88
29. Jimbo, M., Misra, K. C., Miwa, T., Okado, T.: Combinatorics of representations of  $U_q(\widehat{\mathfrak{sl}}(n, \mathbf{C}))$  at  $q = 0$ . Preprint RIMS **709** (1990). *Comm. Math. Phys.* (to appear)
30. Date, E., Jimbo, M., Kuniba, A., Miwa, T., Okada, M.: Paths, Maya diagrams and representations of  $\widehat{\mathfrak{sl}}(2, r)$ . *Adv. Studies Pure Math.* **19** (1989) 149–191
31. Kac, V. G., Wakimoto, M.: Modular and conformal invariance constraints in representation theory of affine Lie algebras. *Adv. Math.* **70** (1988) 156–236
32. Date, E., Jimbo, M., Kuniba, A., Miwa, T., Okada, M.: One-dimensional configuration sums in vertex models and affine Lie algebra characters. *Lett. Math. Phys.* **17** (1989) 69–77

# The Periodic Problems for Two-Dimensional Integrable Systems

Igor Krichever

Landau Institute for Theoretical Physics, Academy of Sciences of the USSR  
GSP-1, 117940 ul. Kosygina 2, Moscow, USSR

## 1. Introduction

Since the middle of the seventies algebraic geometry has become a very powerful tool in various problems of mathematical and theoretical physics. In the theory of integrable equations the algebraic geometrical methods provide a construction of the periodic and quasi-periodic solutions which can be written exactly in terms of theta functions of the auxiliary Riemann surfaces.

All the integrable equations which are considered in the soliton theory can be represented as compatibility conditions of the auxiliary linear problems. One of the most general types of such representations has the form:

$$[\partial_y - L, \partial_t - A] = 0, \quad (1.1)$$

where  $L, A$  are differential operators of the form

$$L = \sum_{i=0}^n u_i(x, y, t) \partial_x^i, \quad A = \sum_{i=0}^m v_i(x, y, t) \partial_x^i \quad (1.2)$$

with scalar or matrix coefficients.

The most important example of these equations is the Kadomtsev-Petviashvili (KP) equation

$$\frac{3}{4} \sigma^2 u_{yy} + \left( u_t - \frac{3}{2} u u_x + \frac{1}{4} u_{xxx} \right)_x = 0 \quad (1.3)$$

which is equivalent to (1.1), where

$$L = \sigma(-\partial_x^2 + u(x, y, t)), \quad A = \partial_x^3 - \frac{3}{2} u \partial_x - w(x, y, t). \quad (1.4)$$

The algebraic geometrical construction of the solutions of integrable equations is based on the concept of the Baker-Akhiezer functions which are defined by their very specific analytical properties on the auxiliary Riemann surfaces. For example, the Baker-Akhiezer functions in the case of the KP equation are defined for each smooth algebraic curve  $\Gamma$  (Riemann surface of finite genus  $g$ ) with the fixed point  $P_0$  on it, and the local parameter  $k^{-1}(P)$  in the neighbourhood of this point,  $k^{-1}(P_0) = 0$ . For any set of generic points  $\gamma_j$ ,  $j = 1, \dots, g$ , there exists a unique function  $\Psi(x, y, t, P)$ ,  $P \in \Gamma$ , such that:

1<sup>0</sup>. It is meromorphic on  $\Gamma$  outside the point  $P_0$  and has no more than simple poles at the points  $\gamma_j$  (if they are distinct);

2<sup>0</sup>. The function  $\Psi$  has the form:

$$\Psi(x, y, t, P) = \left( 1 + \sum_{s=1}^{\infty} \xi_s^i(x, y, t) k^{-1} \right) \exp(ikx + \sigma^{-1}k^2y + ik^3t) \quad (1.5)$$

$k = k(P)$ , near the point  $P_0$ .

For any formal series of the form (1.5) there exist unique operators  $L$  and  $A$  of the form (1.4) such that the following relations

$$\begin{aligned} (\partial_y - L)\Psi &= O(k^{-1}) \exp(ikx + \sigma^{-1}k^2y + ik^3t) \\ (\partial_t - A)\Psi &= O(k^{-1}) \exp(ikx + \sigma^{-1}k^2y + ik^3t) \end{aligned} \quad (1.6)$$

are valid. From (1.6) it follows that the coefficient  $u(x, y, t)$  of these operators is equal to

$$u(x, y, t) = 2i\xi_{1,x}(x, y, t). \quad (1.7)$$

The left hand sides of (1.6) define the functions which have the same analytical properties outside  $P_0$  as  $\Psi$ , and have the form (1.6) near this point. From the uniqueness of the Baker-Akhiezer function  $\Psi$ , it follows that they are equal to zero. Hence,

$$(\partial_y - L)\Psi = 0, \quad (\partial_t - A)\Psi = 0 \quad (1.8)$$

and  $u(x, y, t)$ , which is given by (1.7) is a solution of the KP equation.

The Baker-Akhiezer function  $\Psi(x, y, t, P)$  can be exactly written in terms of the Abelian differentials and Riemann theta-function. From the corresponding formulae it follows that the above constructed solutions of the KP equation have the form

$$u(x, y, t) = 2\partial_x^2 \ln \theta(Ux + Vy + Wt + \Phi/\tau) + \text{const}. \quad (1.9)$$

Here,  $\theta(z_1, \dots, z_g)$  is the Riemann theta-function which is defined by the matrix  $\tau_{ij}$  of the  $b$ -periods of the normalized holomorphic differentials on  $\Gamma$ . The vectors  $2\pi iU, 2\pi iV, 2\pi iW$  are the vectors of  $b$ -periods of the normalized Abelian differentials of the second kind with the only poles at  $P_0$  of orders 2, 3, 4, respectively. The vector  $\Phi$  corresponds to the set of the points  $\gamma_j$  and can be considered in (1.9) as an arbitrary vector.

The construction was proposed in [1, 2] and was developed in different ways for various types of integrable equations (see, for example, the reviews [3, 4, 5, 6]. The analytical properties of the Baker-Akhiezer functions are the natural generalization of the analytical properties of Bloch functions of the ordinary periodic finite-gap differential operators which were obtained in the remarkable works by Novikov, Dubrovin, Matveev, Its in which the algebraic geometrical solutions of the KdV equation, sine-Gordon equation and some other Lax-type equations were constructed.

In this report we shall present a brief review of the latest results obtained in the theory of periodic problems for the two-dimensional integrable systems. First of all, why is it algebraic geometry? What is the meaning of the algebraic geometrical solutions for the general periodic (in  $x$  and  $y$ ) initial value problem for such equations? For the one-dimensional evolution integrable equations, the algebraic geometrical solutions are dense in the space of all periodic solutions

(though this statement has not been proved rigorously for all such equations). In the case of the two-dimensional integrable equations the situation is much more complicated.

There are two real forms of the KP-1 ( $\sigma^2 = -1$ ) and KP-2 ( $\sigma^2 = 1$ ). It turns out, that the periodic problems for these equations differ dramatically from each other.

The formal non-integrability of the periodic problem for the KP-1 equation was proved in [7]. The proof of the integrability of such problem for the KP-2 equation was obtained by the author [8] and is based on the spectral theory of the operator

$$M = \sigma \partial_y - \partial_x^2 + u(x, y) \quad (1.10)$$

with the periodic potential.

The second problem which will be considered in this talk is the perturbation theory for two-dimensional integrable equations. We shall concentrate our attention on the so-called Whitham equation which is in our case a system of equations on bundles over the Teichmüller spaces. Finally, we shall demonstrate how the Whitham theory and other aspects of the perturbation theory of integrable equations will be married to each other in attempts to solve the Heisenberg relations

$$[L_n, A_m] = 1, \quad (1.11)$$

for the ordinary differential linear operators

$$L_n = \sum_{i=0}^n u_i(x) \partial_x^i, \quad A_m = \sum_{i=0}^m v_i(x) \partial_x^i, \quad u_n = v_m = 1. \quad (1.12)$$

The latter are the most popular subject in the field of string theory.

## 2. The Spectral Theory of Two-Dimensional Periodic Linear Differential Operators

The solutions  $\Phi(x, y, w_1, w_2)$  of the nonstationary Schrödinger equation

$$(\sigma \partial_y - \partial_x^2 + u(x, y)) \Phi(x, y, w_1, w_2) = 0 \quad (2.1)$$

with the periodic potential are called Bloch solutions, if they are eigenfunctions of the monodromy operators, i.e.

$$\begin{aligned} \Psi(x + a_1, y, w_1, w_2) &= w_1 \Psi(x, y, w_1, w_2) \\ \Psi(x, y + a_2, w_1, w_2) &= w_2 \Psi(x, y, w_1, w_2). \end{aligned} \quad (2.2)$$

The set of pairs  $Q = (w_1, w_2)$ , for which there exists such a solution is called the Floque set, and will be denoted by  $\Gamma$ . The multivalued functions  $p(Q), E(Q)$  such that

$$w_1 = \exp(ipa_1), \quad w_2 = \exp(iEa_2)$$

are called quasi-momentum and quasi-energy, respectively.

For the “free” operator with zero potential  $u_0 = 0$ , the Floque set is parametrized by the points of the complex  $k$ -plane

$$w_1^0 = \exp(ika_1), \quad w_2^0 = \exp(-\sigma^{-1}k^2a_2) \quad (2.3)$$

and the Bloch solutions have the form

$$\Psi_0(x, y, k) = \exp(ikx - \sigma^{-1}k^2y). \quad (2.4)$$

It turns out that if  $\operatorname{Re} \sigma \neq 0$ , then the Floque set of the operator (2.1) with the smooth potential  $u(x, y)$  is isomorphic to the Riemann surface  $\Gamma$  (which has in a generic case infinite genus). The corresponding Riemann surface has such a specific structure that the theory of abelian differentials, theta-functions and so on can be constructed for it as well as for the finite genus case.

The source of the difference between the two cases  $\operatorname{Re} \sigma = 0$  and  $\operatorname{Re} \sigma \neq 0$  is the difference between the structure of the “resonant” points for the free operators. The resonant points are the points on the complex  $k$ -plane which are the pre-images of the self-intersection points of the imbedding  $C \rightarrow C^2$ , which is defined by (2.3). The points  $k$  and  $k'$  are resonant, if

$$w_i^0(k) = w_i^0(k'), \quad i = 1, 2. \quad (2.5)$$

From (2.3) it follows that such points are parametrized by integers ( $N > 0, M$ ) and have the form:

$$k = k_{N,M}, \quad k' = k_{-N,-M}, \quad (2.6)$$

where

$$k_{N,M} = \frac{\pi N}{a_1} + i\sigma \frac{Ma_1}{Na_2}. \quad (2.7)$$

In case  $\operatorname{Re} \sigma \neq 0$ , the resonant points tend to infinity and, hence, have no limiting points outside infinity. In case  $\operatorname{Re} \sigma = 0$ , the resonant points are dense on the real axis which makes it impossible (at least by means of our methods) to construct the global Riemann surface of the Bloch functions.

For the real smooth potential  $u$  the Floque set can be described in the following form. Let us call the set of pairs of the complex numbers  $\pi = \{p_{s,1}, p_{s,2}\}$  (where  $s$  belongs to any finite or infinite subset of integer pairs ( $N > 0, M$ )) “admissible”, if

$$\operatorname{Re} p_{s,i} = \frac{\pi N}{a_1}, \quad |p_{s,i} - k_s| = O\left(\frac{1}{|k_s|}\right), \quad i = 1, 2$$

and the intervals  $[p_{s,1}, p_{s,2}]$  do not intersect each other. Let us define the Riemann surface  $\Gamma(\pi)$  for any admissible set  $\pi$ . It is obtained from the complex  $k$ -plane by cutting it along the intervals  $[p_{s,1}, p_{s,2}]$  and  $[-\bar{p}_{s,1}, -\bar{p}_{s,2}]$  and by sewing after that the left side of the first cut with the right side of the second cut and vice versa.

**Theorem 1.** *For any real periodic potentials  $u(x, y)$ , which can be analytically extended in some neighbourhood of the real values  $x, y$ , the Bloch solutions of the Equation (2.1) with  $\sigma = 1$  are parametrized by the points  $Q$  of the Riemann surface  $\Gamma(\pi)$  corresponding to some admissible set  $\pi$ . The function  $\Psi(x, y, Q)$  which is normalized by the condition  $\Psi(0, 0, Q) = 1$ , is meromorphic on  $\Gamma$  and has a simple pole  $\gamma_s$  on each cycle  $a_s$  which corresponds to the cut  $[p_{s,1}, p_{s,2}]$ . If the admissible set  $\pi$  contains only a finite number of pairs, then  $\Gamma(\pi)$  has finite genus and is compactified by only one point  $P_0(k = \infty)$ , in the neighbourhood of which the Bloch function  $\Psi$  has the form (1.5).*

The potentials  $u$  for which  $\Gamma(\pi)$  has finite genus, are called finite-gap potentials and as it follows from the last statement of the theorem that they coincide with the algebraic geometrical potentials.

**Theorem 2.** *Any smooth periodic potential  $u$  of the Equation (2.1) (with  $\operatorname{Re} \sigma \neq 0$ ), which can be analytically extended in the complex neighbourhood of the real  $x, y$ , can be approximated uniformly with any number of the derivatives by means of the finite-gap (algebraic geometrical) potentials.*

The Floque set is the “integral” of the KP equation. From the previous theorems we have:

**Theorem 3.** *For any smooth periodic function  $v(x, y)$  there exists a unique solution of the KP-2 equation  $u(x, y, t)$ , such that  $u(x, y, 0) = v(x, y)$ . This solution is regular for all  $t$  and quasi-periodic in  $t$ . Any smooth periodic solutions of the KP-2 equation can be approximated by means of the finite-gap solutions.*

### 3. The Perturbation Theory of the Finite-Gap Solutions. Whithem Equations

The non-linear WKB (or Whithem) method can be applied to any non-linear equation which has the set of the exact solutions of the form

$$u_0(x, y, t) = u_0(Ux + Vy + Wt + \Phi|I_1, \dots, I_N), \quad (3.1)$$

where  $u_0(z_1, \dots, z_g|I_k)$  is a periodic function of the variable  $z_i$  depending on the parameters  $I_k$ . The vectors  $U, V, W$  are also functions of the same parameters:  $U = U(I), V = V(I), W = W(I)$ .

In the framework of the non-linear WKB-method the asymptotic solutions of the form

$$u(x, y, t) = u_0(\varepsilon^{-1} S(X, Y, T)|I_k) + \varepsilon u_1 + \dots \quad (3.2)$$

are constructed for the perturbed or non-perturbed initial equation. Here  $X = \varepsilon x$ ,  $Y = \varepsilon y$ ,  $T = \varepsilon t$  are the “slow variables”. If the vector  $S(X, Y, T)$  is defined from the relations

$$\begin{aligned} \partial_X S &= U(I(X, Y, T)) = U(X, Y, T) \\ \partial_Z S &= V(X, Y, T), \quad \partial_T S = W(X, Y, T) \end{aligned} \quad (3.3)$$

the main term  $u_0$  in the expansion (3.2) satisfies the initial equation up to the first order in  $\varepsilon$ . After that all the other terms of the series (3.2) are defined from the non-homogeneous linear equations. The construction of such asymptotic solutions even for integrable equations is very important, because when using the slow modulation of the parameters of their exact solutions one can sometimes solve the integrable equation with “non-integrable boundary conditions”.

For the KdV equation and for some other Lax-type equations, the Whithem method was developed and applied to various problems in [9, 10, 11]. For the two-dimensional integrable systems the Whithem method was proposed in [12]. We shall present here only a part of the corresponding results.

The asymptotic solutions of the form (3.2) can be constructed with an arbitrary dependence of the parameters  $I_k$  on slow variables. In this case the expansion

(3.2) is valid on the scales of order 1. The right hand side of the non-homogeneous linear equation which defines the first order term  $u_1$  contains the first derivatives of the parameters  $I_k$ . Therefore, the choice of the dependence of  $I_k$  on slow variables can be used for the cancellation of the “secular” term in  $u_1$ . The corresponding equations on  $I_k$  are usually called the Whitham equations.

Let us consider again the KP equation as an example of the two-dimensional integrable systems. Its finite-gap solutions have the form (3.1). The set of their parameters are the system of local coordinates of the manifold  $M_g$  which has dimension  $N = 3g + 1$ .

$$M_g((\Gamma, P_0[k^{-1}]_2) . \quad (3.4)$$

(Two local parameters are  $m$ -equivalent if  $k_1 = k + O(k^{-m})$ ; the corresponding equivalence class of the local parameter is denoted by  $[k^{-1}]_m$ .)

Let us consider the second kind differentials on  $\Gamma$  with the only poles at the point  $P_0$  of the form

$$dp = dk(1 + O(k^{-2})), \quad dE = i\sigma^{-1}dk^2(1 + O(k^{-3})), \quad d\Omega = dk^3(1 + O(k^{-4})) \quad (3.5)$$

which have the real periods for any cycle on  $\Gamma$ . Their integrals  $p(Q)$ ,  $E(Q)$ ,  $\Omega(Q)$  are multivalued functions on the manifold  $M_g^*$  which is a bundle over  $M_g$

$$M_g^* = (\Gamma, P_0, [k^{-1}]_2, Q \in \Gamma) . \quad (3.6)$$

If  $(\lambda, I_1, \dots, I_{3g+1})$  is a system of local coordinates on  $M_g^*$  and  $I_k$  are functions of the variables  $X, Y, T$  then  $p = p(\lambda, X, Y, T)$ ,  $E = E(\lambda, X, Y, T)$ ,  $\Omega = \Omega(\lambda, X, Y, T)$  become functions of these variables.

**Theorem 4.** *The necessary conditions for the existence of the asymptotic solutions of the equation*

$$\frac{3}{4}\sigma^2 u_{yy} + (u_t - \frac{3}{2}uu_x + \frac{1}{4}u_{xxx})_x + \varepsilon K[u] = 0 \quad (3.7)$$

*which has the form (3.2) with uniformly bounded first-order term are equivalent to the equation*

$$\frac{\partial p}{\partial \lambda} \left( \frac{\partial E}{\partial T} - \frac{\partial \Omega}{\partial T} \right) - \frac{\partial E}{\partial \lambda} \left( \frac{\partial p}{\partial T} - \frac{\partial \Omega}{\partial X} \right) + \frac{\partial \Omega}{\partial \lambda} \left( \frac{\partial p}{\partial Y} - \frac{\partial E}{\partial X} \right) = \frac{\langle \Psi K \Psi_+ \rangle_x}{\langle \Psi \Psi_+ \rangle_x} \frac{\partial p}{\partial \lambda}. \quad (3.8)$$

Here  $K[u]$  is an arbitrary differential polynomial;  $\Psi, \Psi_+$  are the corresponding Baker-Akhiezer function and its dual, respectively.

*Remark.* It turns out that there are only  $3g + 1$  independent equations among the Equation (3.8) which should be fulfilled for any point  $Q$  of the curve  $\Gamma$ .

For the KdV equation and  $K = 0$  the Equation (3.8) have the form

$$\partial_T p = \partial_X \Omega \quad (3.9)$$

which was obtained in [11]. The construction of the exact solutions of the Equation (3.8) with  $K = 0$  was proposed in the work [12]. We shall present the particular case of this scheme in the next section where the Heisenberg relations would be considered.

#### 4. The Heisenberg Relations for the Ordinary Linear Differential Operators

Great progress has been made recently in non-perturbative two-dimensional gravity coupled to various matter fields. It was shown that the dependence of physical quantities (such as specific heat) on scaled coefficients of the models is described by the KP-hierarchy on the space of the ordinary linear differential operators  $L_n, A_m$  such that the relations (1.11) are fulfilled. For pure two-dimensional gravity  $n = 2, m = 3$  the Equation (1.11) is equivalent to the Painlevé 1 equation

$$\frac{1}{4}u_{xxx} - \frac{3}{2}uu_x = 1. \quad (4.1)$$

The Equation (1.12) has a simple scaling transformation

$$u_i(x) = \varepsilon^{(i-n)\beta} \tilde{u}_i(\varepsilon^{-\beta}x), \quad v_i = \varepsilon^{(i-m)\beta} \tilde{v}_i(\varepsilon^{-\beta}x) \quad (4.2)$$

$\beta = (n+m)^{-1}$ . For the operators  $\tilde{L}_n, \tilde{A}_m$  with the coefficients  $\tilde{u}_i, \tilde{v}_i$  we have

$$[\tilde{L}_n, \tilde{A}_m] = \varepsilon. \quad (4.3)$$

The formal asymptotic solutions of the equation (4.3) can be constructed using any commuting operators  $[L_{n,0}, A_{m,0}] = 0$

$$\tilde{L}_n = L_{n,0} + \varepsilon L_{n,1} + \dots, \quad \tilde{A}_m = A_{m,0} + \varepsilon A_{m,1} + \dots. \quad (4.4)$$

Unfortunately, these asymptotic solutions are well-defined only in the interval  $x \sim 1$ . For our purposes it is necessary to have the solutions for  $x \sim \varepsilon^{-1/(n+m)}$ . It can be done in framework of the Whitham theory.

The commuting operators of co-prime orders  $(n, m) = 1$  are parametrized by the coefficients of the polynomial

$$w^n + E^m + \sum_{in+jm \leq nm-2} \alpha_{ij} w^i E^j = 0 \quad (4.5)$$

and by the points of the Jacobian of the corresponding algebraic curve [13]. In [1, 2] the exact formulae for the coefficients of the generic commuting operators in terms of the Riemann theta-function were found. For example,

$$u_{n-2} = -n\partial_x^2 \ln \theta(Ux + \Phi/\tau) + \text{const}. \quad (4.6)$$

Here the matrix  $\tau$  of  $b$ -periods of  $\Gamma$  depends on the values  $\alpha_{ij}$ . The vector  $U$  is also a function of the variables  $\alpha_{ij}$ . The phase vector  $\Phi$  is arbitrary. All the other coefficients have the same structure

$$u_i = u_{i,0}(Ux + \Phi/\tau), \quad v_i = v_{i,0}(Ux + \Phi/\tau). \quad (4.7)$$

Let me consider the operators  $L_n^\#, A_m^\#$  with the coefficients

$$u_i^\# = u_{i,0} \left( \frac{1}{\varepsilon} S(X)/\tau(X) \right), \quad v_i^\# = v_{i,0} \left( \frac{1}{\varepsilon} S(X)/\tau(X) \right). \quad (4.8)$$

If the vector  $S(X)$  is defined by the relation  $\partial_X S = U(\alpha_{ij}(X))$  then the operators  $L_n^\#, A_m^\#$  commute up to the order  $\varepsilon$ . As was shown in [12] in more general situation the requirement that the first order terms in the expansion (4.4) should be uniformly bounded for all  $x$  leads to the equations on the variables  $\alpha_{ij}$ . They are particular cases of (3.8) and have the form

$$\frac{\partial w(E, X)}{\partial X} = \frac{\partial p(E, X)}{\partial E}. \quad (4.9)$$

It turns out that they are integrable and we present the construction of their solutions below. Our conjecture (which is partly proved now for  $n = 2, m = 3$ ) is that all the other terms of the asymptotic solutions (4.4) are also bounded and the series (4.4) are convergent. If this is true, it is possible to make the inverse rescaling and find the limit for  $\varepsilon \rightarrow 0$ . To begin with we shall give the final answer for the KdV equation with the “string” boundary conditions (1.12)  $n = 2$ ,  $m = 2k + 1$ .

Let us consider an arbitrary hyperelliptic curve  $\Gamma$

$$y^2 = \prod_{i=1}^{2k+1} (E - E_i) = E^{2k+1} + \prod_{i=1}^{2k} c_i E^i = R(E). \quad (4.10)$$

As is well-known, this curve defines the solutions of the KdV equation which have the form (1.9) (with  $V = 0$ ).

For any given set of the parameters: the complex constants  $c_{k,0}, c_{k+1,0}, \dots, c_{2k,0}$ , the real constants  $h_i, h'_i, i = 1, \dots, k$ , we shall consider the hyperelliptic curve which is defined by the polynomial  $R$  with the coefficients

$$c_i = c_{i,0}, \quad i = k + 2, \dots, 2k; \quad c_k = x + c_{k,0}; \quad c_{k+1} = t + c_{k+1,0} \quad (4.11)$$

and such that

$$\operatorname{Im} \int_{E_{2i}}^{E_{2i+1}} \sqrt{R} dE = h_i, \quad \operatorname{Im} \int_{E_1}^{E_{2i}} \sqrt{R} dE = h'_i, \quad i = 1, \dots, k. \quad (4.12)$$

The Equation (4.12) are the set of  $2k$  real equations which define  $k$  unknown complex coefficients  $c_i, i = 0, \dots, k - 1$  of the polynomial  $R(E)$ . They become functions of the variable  $x, t$ . The  $\tau$  matrix of the corresponding curve becomes a (known) function of the variables  $x, t$ . Let us define the vector

$$S_i(x, t) = \frac{1}{\pi} \left( \int_{E_1}^{E_{2i}} \sqrt{R} dE - \sum_{j=1}^k \tau_{ij} \int_{E_{2j}}^{E_{2j+1}} \sqrt{R} dE \right). \quad (4.13)$$

**The Main Conjecture.** *The functions*

$$u(x, t) = -2\partial_x^2 \ln \theta(S(x, t) + \Phi/\tau(x, t)) - 2r_1(x, t) \quad (4.14)$$

*are the exact solutions of the KdV equation with the “boundary conditions” (1.12).*

Here  $r_1(x, t)$  is the coefficient of the differential

$$d\Omega_1 = \frac{E^k + \sum_{i=0}^{k-1} r_i E^i}{2\sqrt{R}} dE, \quad \int_{E_{2l}}^{E_{2l+1}} d\Omega_1 = 0.$$

Thus the Equations (4.12) are the only transcendental equations in the definition of  $u(x, t)$ .

Let us consider now the general Heisenberg relations. Any equation of the form (4.5) has the formal solution

$$w = k^m + \sum_{i=-m+2}^{\infty} a_i k^{-i}, \quad k^n = E. \quad (4.15)$$

This means that the affine curve (4.5) is compactified by a single point  $P_0$ . Let us fix a first few coefficients of the expansion (4.15) and denote them by

$$a_{n-j} = \frac{j}{n} t_j, \quad j = 1, \dots, m+n-2. \quad (4.16)$$

They uniquely define the following coefficients of (4.5)

$$\alpha_{ij}, im + jn \geq (m-1)(n-1) = 2g. \quad (4.17)$$

For any given real numbers  $h_i, h'_i$ ,  $i = 1, \dots, g$ , all the other coefficients of the polynomial (4.5) can be defined (at least locally) as functions of the parameters  $t_j$  with the help of the relations

$$\operatorname{Im} \int_{a_i} w dE = h_i, \quad \operatorname{Im} \int_{b_i} w dE = h'_i. \quad (4.18)$$

They give  $2g$  real equations on  $g$  complex variables  $\alpha_{ij}, in + jm < (n-1)(m-1)$ . Therefore, the curve  $\Gamma$  and the algebraic function  $w(E)$  become functions of the variables  $t_j$ .

**Theorem 5.** *The function  $w(E, t_1, \dots)$  satisfies the Whitham equations (4.9) if  $t_1 = x$ .*

Let us define the differentials  $d\Omega_j$ ,  $j = 1, \dots, m+n-2$ , whose only poles at infinity have the form

$$d\Omega_j = dk^j (1 + O(k^{-j-1})) \quad (4.19)$$

and such that

$$\operatorname{Im} \int_{\gamma} d\Omega_j = 0, \quad \gamma \in H_1(\Gamma). \quad (4.20)$$

**Corollary.** *If the relations (4.18) are fulfilled, then*

$$\frac{\partial p}{\partial t_j} = \frac{\partial \Omega_j}{\partial x} \quad \frac{\partial \Omega_j}{\partial t_j} = \frac{\partial \Omega_j}{\partial t_i}.$$

*Remark.* It can be shown that the conjecture which was proposed recently in [14] leads to one particular solution of the Painlevé 1 which belongs to our set of solutions.

## References

1. Krichever, I.M.: An algebraic-geometrical construction of the Zakharov-Shabat equation and their periodic solutions. *Dokl. Akad. Nauk SSSR* **227**, 291–294
2. Krichever, I.M.: The integration of non-linear equations by methods of algebraic-geometry. *Funkt. Anal. Priloz.* **11** (1) (1977) 15–31
3. Krichever, I.M.: Methods of algebraic geometry in the theory of non-linear equations. *Uspekhi Matem. Nauk* **32** (6) (1977) 183–208
4. Dubrovin, B.A.: Theta-functions and non-linear equations. *Uspekhi Matem. Nauk* **36** (6) (1981) 11–80
5. Krichever, I.M., Novikov, S.P.: Holomorphic bundles over algebraic curves and non-linear equations. *Uspekhi Matem. Nauk* **35** (6) (1980) 47–68
6. Dubrovin, B.A., Krichever, I.M., Novikov, S.P.: Integrable systems. Dynamical systems **4**, Itogi Nauki i Tekhniki, Fund. invest., VINITI Akad. Nauk SSSR, 1985
7. Zakharov, V.E., Schulmann, E.I.: On problems of the integrability of two-dimensional systems. *Dokl. Akad. Nauk SSSR* **283** (6) (1985) 1325–1329
8. Krichever, I.M.: Spectral theory of two-dimensional periodic operators and its applications. *Uspekhi Matem. Nauk* **44** (2) (1989) 121–184
9. Gurevich, A.V., Pitaevskii, L.P.: Nonstationary structure of the interactionless shock-wave. *JETP* **65** (3) (1973) 590–604
10. Dobrokhotov, S.Yu., Maslov, V.P.: Multiphase asymptotics of non-linear partial differential equations with a small parameter. Soviet Scientific Reviews, Math. Phys. Rev. **3** (1982) 221–280, OPA Amsterdam
11. Flashka, H., Forest, M., McLaughlin, L.: The multiphase averaging and the inverse spectral solution of the Korteweg-de Vries equation. *Comm. Pure Appl. Math.* **33** (6) (1980) 739–784
12. Krichever, I.M.: The averaging method for two-dimensional integrable equations. *Funkt. Anal. Priloz.* **22** (3) (1988) 37–52
13. Burchnall, J.L., Chaundy, T.W.: Commuting ordinary differential operators, II. *Proc. Roy. Soc. London* **118** (1928) 557–583
14. Novikov, S.P.: On relations  $[L, A] = 1$ . To appear in *Funkt. Anal. Priloz.* 1990

# Renormalization Group and Random Systems

*Antti Kupiainen*

Rutgers University, Department of Mathematics, New Brunswick, NJ 08903, USA

## 1. Introduction

Some of the most interesting and challenging problems in theoretical and mathematical physics have been the ones involving in an essential way several distance scales. Such problems include the theory of critical phenomena and phase transitions, the problem of the existence of relativistic quantum field theories, the turbulent behaviour in hydrodynamic systems and many others. While many of these problems date back several decades (or, indeed, to the previous century), significant progress was made only in the 60s, mainly in the context of critical phenomena, using ideas stressing the role of scale invariance in such problems. This line of thought culminated with the creation of the theory of the Renormalization Group (RG) by K.Wilson (see [1] for the history) that provided a unified method of attack for multiscale problems. The 70s saw the RG-theory vindicated in approximative analytic treatments of a multitude of problems in critical phenomena, hydrodynamics and other many-body systems. With the increasing power of computers, the ideas entered numerical analysis of such problems in the 80s.

On the mathematical side, RG ideas lead to the rigorous theory of so-called renormalizable quantum field theories (and a non-renormalizable one too) and some critical statistical mechanics problems [2] (these works by no means exhaust the mathematical progress in these fields, nor even the RG inspired works, for a rather recent set of references, see [3]. In this talk I would like to review some recent attempts to develop a mathematical RG theory for the study of so-called random (or disordered) systems. Again, this field has seen plenty of beautiful mathematical work, some of it RG inspired too, but no attempt is being made here for a comprehensive review.

## 2. Disordered Spin Systems

The mathematical framework of the problems we are going to study is that of *path space measures*. This is a unified framework for classical statistical mechanics, Euclidean quantum field theory, diffusion processes and many other systems. One

considers measures on some space of maps

$$\phi : \mathcal{S} \rightarrow \mathcal{M}$$

where  $\mathcal{S}$  is the “space” or “time” or “space-time” and  $\mathcal{M}$  the “field-space”. For QFT,  $\mathcal{S}$  is a manifold ( $\mathbf{R}^n$ , a Riemann surface...), for statistical mechanics a lattice (e.g.  $\mathbf{Z}^n$ ) and  $\mathcal{M}$  can be a manifold, a group or just some set, depending on the concrete model.

We will concentrate in this talk mainly on two examples of the above, which occur in a wide variety of physical contexts. The first is the *Ising Model*, where  $\mathcal{S}$  is  $\mathbf{Z}^d$  and  $\mathcal{M}$  is  $\mathbf{Z}_2 = \{1, -1\}$ . To define the measure on the set of such maps  $\mathcal{F} = (\mathbf{Z}_2)^{\mathbf{Z}^d}$ , one first considers a subset of  $\mathcal{F}$ ,  $\mathcal{F}_A^\pm = \{\phi \mid \phi(x) = \pm 1, x \in A^c\}$  and on this finite set (take  $A$  finite) the probability measure

$$\mu_A^{\beta, \pm}(\{\phi\}) = Z^{-1} e^{\beta \sum (\phi(x)\phi(y)-1)} \quad (1)$$

where  $x$  and  $y$  run through nearest neighbors on the lattice. The measures  $\mu^{\beta, \pm}$  are constructed as limits of (1) as  $A \rightarrow \mathbf{Z}^d$ . They describe *symmetry breaking*: for  $\beta$  large and  $d > 1$ ,

$$\lim_{A \rightarrow \mathbf{Z}^d} \int \phi(x) d\mu_A^{\beta, \pm}(\phi) = \int \phi(x) d\mu^{\beta, \pm}(\phi) = \pm m \quad (2)$$

with  $m \neq 0$ . The boundary conditions  $\pm$  in (1) select non-zero *magnetization*  $\pm m$  at the infinite volume limit.

Let us now present a disordered version of (1). The physical idea behind the disorder is that real materials, which are modelled by (1), have impurities, such as atoms of different kind randomly placed in the crystal structure. Mathematically, this translates into spatially varying randomness in the various parameters of the model. It turns out that a large class of disorder is modelled by replacing the exponent in (1) by

$$\beta \left( \sum_{x,y \in \mathbf{Z}^d} J(x,y)(\phi(x)\phi(y)-1) + \sum_{x \in \mathbf{Z}^d} (\phi(x) - (\pm 1))h(x) \right). \quad (3)$$

The functions  $J$  and  $h$  are *random* with some prescribed probability distribution. Thus, as  $A \rightarrow \mathbf{Z}^d$ , we get *random measures*  $\mu^{J,h}(\phi)$  on  $\mathcal{F}$  and can ask whether the  $m$  in (2) is nonzero almost everywhere in  $J$  and  $h$  (this is, of course, only interesting if the randomness preserves the  $\phi \rightarrow -\phi$  symmetry of the system).

An example of (3) is the *Spin Glass*, where  $h = 0$  and e.g.  $J(x,y) = 0$ , unless  $|x-y|=1$  and the non-zero  $J$ 's are taken independent and identically distributed, with a  $J \rightarrow -J$  symmetric distribution.

The example we will discuss below is the *Random Field Ising Model*, given by (3) with  $J$  non-random as in (1) and the  $h(x)$ ,  $x \in \mathbf{Z}^d$  independent, identically distributed, with a distribution invariant under  $h(x) \rightarrow -h(x)$  and with variance  $Eh(x)^2 = \epsilon^2$ .

**Theorem 1** [4, 5, 6]. Let  $d > 2$ . Then, there exists a  $\beta_0 < \infty$  and an  $\varepsilon_0 > 0$  such that for  $\beta > \beta_0, \varepsilon < \varepsilon_0$  the limit (2) exists and  $m \neq 0$  with probability 1. For  $d \leq 2$  the magnetization  $m$  vanishes for all  $\beta$  and  $\varepsilon$ .

**Remark.** Let us define the *lower critical dimension*  $d_L$  as the largest  $d$  where the magnetization vanishes for all  $\beta$ . Then, the Theorem states that  $d_L = 2$ . The corresponding dimension for the non-random system is 1. In the physics literature there was a long controversy on whether  $d_L = 2$  or  $d_L = 3$ . Original arguments due to Imry and Ma [7] suggested that  $d_L = 2$ . These were later challenged by arguments originating in quantum field theory, coined as “dimensional reduction”. This was a rule, partially justified, stating that  $d_L$  for the random system is two more than the one of the non-random one. Since the latter is 1, it was predicted that  $d_L = 3$ . The controversy was solved, at  $\beta = \infty$ , i.e. for the ground state of the Hamiltonian (3), in [4], and for large  $\beta$  in [5]. The  $d = 2$  result is due to [6].

We will outline below the RG proof of the  $d > 2$  part of Theorem 1.

### 3. Diffusion in Random Media

Ordinary diffusion provides another example, where randomness brings interesting effects. Consider random walk on the lattice  $\mathbf{Z}^d$  described by *transition probabilities*  $p(x, y)$  from  $x \in \mathbf{Z}^d$  to  $y \in \mathbf{Z}^d$ :

$$p : \mathbf{Z}^d \times \mathbf{Z}^d \rightarrow [0, 1] \quad (4)$$

satisfying

$$\sum_{y \in \mathbf{Z}^d} p(x, y) = 1. \quad (5)$$

$p$  allows us to define measures  $\mu_T$ ,  $T \in \mathbf{N}$  on the space  $\Omega_T$  of walks  $\omega : \{0, 1, \dots, T\} \rightarrow \mathbf{Z}^d$  starting from  $\omega(0) = 0$ :

$$\mu_T(\{\omega\}) = \prod_{i=1}^T p(\omega(i-1), \omega(i)). \quad (6)$$

Diffusion is a property of the large  $T$  limit of such measures. It will be convenient to realize them as measures  $\nu_T$  on  $C([0, 1])$ , the space of continuous paths  $\omega : [0, 1] \rightarrow \mathbf{R}^d$ , by a simple rescaling. Thus, given an  $\omega \in \Omega_T$ , we obtain a piecewise linear path

$$\tilde{\omega}(t) = T^{-1/2}(\omega(i-1) + (Tt - i + 1)(\omega(i) - \omega(i-1))) \quad (7)$$

where  $i-1 = [Tt]$  and  $[ ]$  denotes the integral part.  $\nu_T$  is the measure induced by (7) on  $C([0, 1])$  with its standard  $\sigma$ -algebra, and we will study the limit

$$\lim_{T \rightarrow \infty} v_T \quad (8)$$

also called the *scaling limit*, and its properties.

The walk is diffusive, if the *diffusion constant*

$$D(p) = \lim_{T \rightarrow \infty} D(T, p) \equiv \lim_{T \rightarrow \infty} T^{-1} \sum_{\omega} \omega(T)^2 \mu_T(\{\omega\}) \quad (9)$$

exists and is non-zero. In terms of the scaling limit

$$D(p) = \int d\nu(\omega) \omega(1)^2. \quad (10)$$

For random walks in *homogenous* environments, the matrix  $p$  is translationally invariant  $p(x, y) = p(x - y)$ . If  $p$  has, say, exponential falloff, the scaling limit is given by the Wiener process. In the disordered system  $p$  is taken as a *random matrix* from some ensemble  $\mathcal{P}$ . One then asks whether the scaling limit and the diffusion constant exist for almost all  $p$ .

Let us now describe the ensemble  $\mathcal{P}$ . We take  $p$  a small random perturbation of the simple random walk:

$$p(x, y) = \frac{1}{2d} + b(x, y) \quad (11)$$

for  $|x - y| = 1$  and  $p(x, y) = 0$  otherwise. Here  $b$  is taken as a random matrix such that

- a)  $b(x, \cdot)$  and  $b(x', \cdot)$  are independent if  $x \neq x'$ , and identically distributed.
- b) The distribution of  $b(x, \cdot)$  is invariant under the natural action of the subgroup of  $O(d)$  fixing the lattice  $\mathbb{Z}^d$ .
- c) The generating function of  $b$  satisfies the bound

$$Ee^{tb(x,y)} < e^{\varepsilon^2 t^2}$$

where  $\varepsilon$  is taken small. So, in particular, the variance of  $b$  is small.

- d)  $\text{Prob}(p(x, y) < e^{-N}) < e^{-\Gamma N}$ , for  $N > 1$ .  $\Gamma$  is taken large.

Under these conditions we have the

**Theorem 2** [8]. *Let  $d > 2$ . Then there exist  $\varepsilon_0 > 0$  and  $\Gamma_0 < \infty$  such that whenever  $\varepsilon < \varepsilon_0$  and  $\Gamma > \Gamma_0$  the diffusion constant  $D(p)$  exists and takes a constant value  $D \neq 0$  a.s. in  $\mathcal{P}$ . The measures  $v_T$  converge weakly to the Wiener measure with diffusion constant  $D$ , a.s. in  $\mathcal{P}$ .*

**Remarks.** 1. Theorem 2 holds for a much wider class of  $p$ 's with exponential falloff in  $|x - y|$  [8]. It also holds for a continuous time version of the problem, where the transition probability  $P(x, t)$  in time  $t$  from 0 to  $x$  satisfies the equation

$$\partial_t P = \Delta P + \nabla \cdot (\mathbf{b}P) \quad (12)$$

where  $\mathbf{b} : \mathbb{Z}^d \rightarrow \mathbf{R}^d$  is random [8].

2. We expect Theorem 2 also to hold for  $d = 2$ . If  $d = 1$ , the situation is completely different [9]. The walk is subdiffusive ; indeed, the mean square distance is proportional to  $(\log t)^4$  with probability tending to 1 as  $t \rightarrow \infty$ . On the other hand, if  $b$  in (11) is taken symmetric i.e.  $b(x, y) = b(y, x)$ , or, if the vector field  $\mathbf{b}$  in (12) is a gradient (of white noise), then under quite general conditions the walk is diffusive in all  $d$  [10].

#### 4. The Renormalization Group

We sketch now the method of proof of Theorems 1 and 2 using the RG. Originally, the RG was developed for the study of (1) (and other similar measures) in a neighbourhood of the *critical point*  $\beta = \beta_c$  at which  $\mu^{\beta, \pm}$  has *long range correlations* (at  $\beta = \beta_c$  the magnetization  $m$  vanishes):

$$\int \phi(x)\phi(y)d\mu \sim C |x - y|^{-\alpha} \quad . \quad (13)$$

as  $|x - y| \rightarrow \infty$ . Here  $\alpha > 0$ . Such long-range correlations make the analysis hard by conventional methods, such as the Taylor expansion in  $\beta$  at  $\beta = 0$  or at  $\beta = \infty$ , since  $\beta_c$  is a point of nonanalyticity of the correlation functions. In a nutshell, the RG approach as applied to such critical situations consists of three steps:

**Coarse Graining.** Choose an integer  $L$  and smear  $\phi$  suitably on scale  $L$  to get a  $\phi_L$ , defined on  $L\mathbb{Z}^n$ .

**Scaling.** Set  $\phi'(x) = L^\gamma \phi_L(Lx)$ . Then  $\phi \rightarrow \phi' \equiv r_L \phi$  defines a map in  $\mathcal{F}$  (often one needs to enlarge  $\mathcal{F}$ , i.e. allow for more general random fields than the  $\phi$  one started with), which induces a map for the measures,  $\mu \rightarrow \mu' \equiv \mathcal{R}_L \mu$ .  $\mathcal{R}_L$  is the RG map.

**Iteration.** For  $L$  fixed, the control of  $\mathcal{R}_L$  is a problem involving short-range correlations and thus often manageable. The solution of the full problem is now translated to the iteration of  $\mathcal{R}_L$ : if  $A$  is, say, a cube of side  $L^N$  in  $\mathbb{Z}^n$ , the integral (2) is done after  $N$  iteration steps. The scale-invariance of the problem is in this approach seen as stabilization of  $\mathcal{R}_L^k \mu$ : this is expected to tend to a *fixed point* of the map  $\mathcal{R}_L$  in a space of probability measures on  $\mathcal{F}$ , provided the scale-parameter  $\gamma$  is chosen properly.

The disordered systems we are considering are not critical in the sense of (13). However, as we will see, they are problems involving several distance scales and the RG approach turns out to be the natural one here too. We first illustrate these steps for the random walks, where the RG will be especially simple.

For coarse graining, given an  $\omega \in \Omega_T$ , let  $\omega_L$  be  $\omega$  restricted to  $L\mathbb{N}$ . Then, for scaling, set

$$\omega'(t) = L^{-1/2} \omega_L(Lt) \equiv (r_L \omega)(t). \quad (14)$$

Note that  $\omega'$  takes values in  $L^{-1/2}\mathbf{Z}^d$ . We get the new measure

$$\mu'(\{\omega'\}) = \mu(\{\omega \mid r_L \omega = \omega'\}) \quad (15)$$

which is given in terms of new transition probabilities

$$\mu'(\{\omega'\}) = \prod_{i=1}^{L^{-1}T} L^{-\frac{d}{2}} p'(\omega'(i-1), \omega'(i)) \quad (16)$$

with

$$p'(x, y) = L^{-d/2} p^L(L^{1/2}x, L^{1/2}y) \quad (17)$$

where  $p^L$  is the  $L$ -th power of the matrix  $p$ . Note that, since  $\omega(i)$  are summed over  $L^{-1/2}\mathbf{Z}^d$ , the powers of  $L$  in (16) are natural. Thus the renormalized measure  $\mu'$  is described in terms of the new matrix  $p'$ , obtained from  $p$  via the non-linear map (17), which we shall call the RG map  $\mathcal{R}$ . In particular we have the identity for the diffusion constants

$$D(t, p) = D(L^{-1}t, \mathcal{R}p) = D(1, \mathcal{R}^n p) \quad (18)$$

if  $t = L^n$ , which means that we have translated the study of the long time behaviour of the walks to the iteration of the RG map.

In (18),  $\mathcal{R}^n p$  are transition probability densities for a walk on  $L^{-n/2}\mathbf{Z}^d$ , and  $\mathcal{R}$  maps such densities to ones for a walk on the finer lattice  $L^{-(n+1)/2}\mathbf{Z}^d$ . In the limit  $n \rightarrow \infty$ , we have walks on  $\mathbf{R}^d$ , and then  $\mathcal{R}$  has a 1-parameter family of Gaussian fixed points

$$p_D^*(x, y) = (2\pi \frac{D}{d})^{-d/2} e^{-d|x-y|^2/2D}. \quad (19)$$

i.e. the transition probabilities of the Wiener process with diffusion constant  $D$ . For example, an exponentially decaying homogenous  $p$  is driven to this fixed point upon iteration of  $\mathcal{R}$ , since in that case  $\mathcal{R}$  is just convolution composed with scaling:

$$\widehat{R^n p}(k) = \hat{p}(L^{-n/2}k)^{L^n} \rightarrow e^{-Dk^2/2d}$$

as  $L \rightarrow \infty$ , where  $\hat{p}$  is the Fourier transform and  $D$  is determined by  $\hat{p}(k) = 1 - Dk^2/2d + \mathcal{O}(|k|^4)$ . This is of course nothing but a reformulation of the central limit theorem.

The main part of the proof of Theorem 2 consists of showing that  $\mathcal{R}^n p \rightarrow p_D^*$  for  $\mathcal{P}$ -almost all  $p$  in a sufficiently strong sense as  $n \rightarrow \infty$ . There are two aspects in this convergence. First, the derivative  $(D\mathcal{R})_p$  turns out to be contractive if  $d > 2$ , with largest eigenvalue  $L^{(2-d)/2}$ . Hence, it turns out that writing  $\mathcal{R}^n p = p_n = Ep_n + b_n$ ,  $b_n$  has variance proportional to

$$\varepsilon_n^2 = L^{(2-d)/2} \varepsilon^2 . \quad (20)$$

The second aspect concerns the condition  $d$ ) above. This is a condition for “traps” in the “environment”: it turns out that, if the bound is violated for  $\Gamma$  small enough, the walk is likely not to be diffusive. The renormalized  $p$  should also satisfy a similar bound (which describes the “traps” in longer distance scales), and indeed, one finds

$$\text{Prob} \left( \int_{\square^c} p_n(x, y) dy < e^{-N} \right) < L^{-n\Gamma} e^{-\Gamma N} \quad (21)$$

where  $\square$  is a unit cube centered at  $x$ . Thus the iteration of the RG “wipes out” the randomness in the matrix  $p$ .

It should be mentioned that, if one starts with  $p$  having correlations with sufficiently slow decay in  $|x - y|$ , the bound (21) may be violated for some  $n$  and thus traps may occur in large scales. Indeed, it is possible to construct  $\mathcal{P}$  where  $p$  have correlations with exponential decay (but slow enough), such that the walk is subdiffusive in all dimensions [11, 12].

Now we turn to the Ising model. It is convenient to represent  $\phi \in \mathcal{F}$  as a (in general disconnected) closed surface  $S$  in the dual lattice: an  $(d-1)$ -cell  $c_{xy}$  dual to a bond  $\{xy\}$  with  $|x - y| = 1$  is in  $S$  if  $\phi(x) \neq \phi(y)$ . Thus, to  $\phi$  there corresponds a surface  $S$  and an assignement of signs  $\pm$  to the components of the complement of  $S$ . The measure  $\mu^{\beta, h}$  (we drop the  $\pm$ ) corresponding to (3) can then be viewed as a measure on the set  $\mathcal{S}$  of all such surfaces

$$\mu^{\beta, h}(\{S\}) = \mathcal{Z}^{-1} e^{-\beta(\text{Area}(S) + \sum_{V_+} h - \sum_{V_-} h)} \quad (22)$$

where  $V_\pm$  are the  $\pm$  regions determined by  $S$ .

If  $h = 0$  and  $\beta$  is large enough, the  $S$ 's are suppressed in the typical configurations of (22). Thus, the probability that the point  $x$  in (2) is in a different component of the complement of  $S$  than  $\infty$  is small and the boundary condition determines the sign at  $x$ . If  $h$  is non-zero, let us consider the case of a connected  $S$ . Then, e.g. for + boundary conditions,

$$\frac{\mu(\{S\})}{\mu(\{\emptyset\})} = e^{-\beta(\text{Area}(S) - 2 \sum_{x \in V_S} h(x))}$$

with  $V_S$  the interior of  $S$ . Such an  $S$  will be probable if

$$h_S = 2 \sum_{x \in V_S} h(x) > \text{Area}(S) . \quad (23)$$

$h_S$ , being a sum of independent random variables, has variance  $Eh_S^2 = |V_S|\varepsilon^2$ , and thus (23) is unprobable if  $d > 2$  since  $\varepsilon|V_S|^{1/2} < \text{Area}(S)$  for all  $S$ . However, we see that disordering configurations, where (23) holds, occur in all scales of the problem. We translate now these observations to the RG language.

For the coarse graining, let  $S \in \mathcal{S}$  be  $S = \cup_i S_i$  with  $S_i$  connected. We set

$$S_L = [\bigcup_{d(S_i) \geq L} S_i] \quad (24)$$

where  $d(S)$  is the diameter of the set and  $[S]$  denotes the smallest union of  $L$ -sided cubes, centered at  $L\mathbb{Z}^d$ , and covering  $S$ .

For scaling, set  $S' = r_L S = L^{-1} S_L$ . Note that  $S'$  is not any more in  $\mathcal{S}$ ; it is a “thick” surface. We set again

$$\mu'(\{S'\}) = \mu^{\beta,h}(\{S \mid r_L S = S'\}). \quad (25)$$

It turns out that  $\mu'$  is approximately of the form (22) with, however, a new  $h'$  having variance  $\sim L^{2-d}$  times the variance of  $h$ , and a new  $\beta' = L^{d-1}\beta$ . These powers are easy to understand: the  $L^{d-1}$  comes from the scaling of the area of surfaces of dimension  $d-1$ , whereas the  $L^{2-d}$  comes from the fact that

$$h'(x) \sim L^{1-d} \sum h(x)$$

with the sum having  $L^d$  independent random variables. The  $L^{1-d}$  is due to the  $\beta$  in (22) multiplying  $h$ .

Thus, approximatively

$$\mathcal{R}^n \mu^{\beta,h} \sim \mu^{\beta_n, h_n} \quad (26)$$

with  $\beta_n$  and the variance  $\varepsilon_n^2$  of  $h_n$  given by

$$\beta_n = L^{n(d-1)} \beta, \quad \varepsilon_n^2 = L^{n(2-d)} \varepsilon^2. \quad (27)$$

The proof of Theorem 1 thus consists in showing that, upon renormalizing, the measure (22) is driven to the trivial  $\beta = \infty$  fixed point with no randomness.

In both cases discussed above, the reason why the disorder gets weaker under  $\mathcal{R}$  is simple: the linearized RG is contractive at the fixed point. However, it is contractive only in the probabilistic sense: e.g. the variance is contracted. There are regions in the space of disorder (in  $\mathcal{P}$  or in the space of  $h$ 's) where  $\mathcal{R}$  is expanding and the main mathematical problem is to show that such regions become more and more unprobable upon the iteration of  $\mathcal{R}$ . These regions correspond to the “traps” in long distance scales or configurations of  $h$  that dominate in (22) the area term.

## 5. Conclusions

There are several problems in the theory of random systems where the RG approach is likely to be the natural one. Among such problems are the classical motion of a particle in the presence of randomly located scatterers (the Lorentz gas), the problem of extended states and diffusion for a Schrödinger operator with random potential and the Spin Glass mentioned above.

For the first two problems it is conceivable that the methods used in [8] for the random walks are sufficient to prove the existence of diffusive behaviour.

In the case of the Spin Glass, there presumably is a non-trivial  $\beta = \infty$  fixed point for the RG. However, here even setting up a useful RG-scheme constitutes a challenge.

## References

1. K. Wilson: The renormalization group and critical phenomena. *Rev. Mod. Phys.* **55** (1983) 583–600
2. K. Gawędzki, A. Kupiainen: Massless lattice  $\phi_4^4$  theory: Rigorous control of a renormalizable asymptotically free model. *Commun. Math. Phys.* **99** (1985) 197–252
  - Gross-Neveu model through convergent perturbation expansions. *Commun. Math. Phys.* **102** (1985) 1–30
  - Renormalization of a non-renormalizable quantum field theory. *Nucl. Phys.* **B262** (1985) 33–48
  - J. Feldman, J. Magnen, V. Rivasseau, R. Seneor: A renormalizable field theory: the massive Gross-Neveu model in two dimensions. *Commun. Math. Phys.* **103** (1986) 67–103
3. J. Glimm, A. Jaffe: Quantum physics. A functional integral point of view. Springer, New York 1987
4. J. Imbrie: The ground state of the three dimensional random field Ising model. *Commun. Math. Phys.* **98** (1985) 145–176
5. J. Bricmont, A. Kupiainen: Phase transition in the 3d random field Ising model. *Commun. Math. Phys.* **116** (1988) 539–572
6. M. Aizenmann, J. Wehr: Rounding effects of quenched randomness on first order phase transitions. *Commun. Math. Phys.*, to appear (1990)
7. Y. Imry, S. K. Ma: Random field instability of the ordered state of continuous symmetry. *Phys. Rev. Lett.* **35** (1975) 1399–1401
8. J. Bricmont, A. Kupiainen: Random walks in asymmetric random environments. Preprint, Rutgers University 1990
9. Y. G. Sinai: *Theor. Prob. Appl.* **27** (1982) 256
10. G. Papanicolaou, S. R. S. Varadhan: In: J. Fritz, J. Lebowitz, D. Szasz (eds.) *Random fields*. North-Holland, 1981, p. 835
11. R. Durrett: *Commun. Math. Phys.* **104** (1986) 87
12. M. Bramson, R. Durrett: *Commun. Math. Phys.* **119** (1988) 119



# Invariants of Links and 3-Manifolds Related to Quantum Groups

Nicolai Reshetikhin \*

Department of Mathematics, Harvard University, Cambridge, MA 02138, USA

Recent years were signified by “strong interaction” between various ideas coming from physics and mathematics. Operator algebra, representation theory, low-dimensional topology,  $q$ -analysis and others are some of the fields involved in this interesting area. Algebraic constructions distilled from the theory of quantum integrable systems are quantum groups.

In the theory of quantum integrable systems it was found in [Ba, Ya] that a certain equation known today as the Yang-Baxter equation is playing the fundamental role in integrability. The essential role of it becomes clear after works by Zamolodchikov's [Z] concerned with factorizable scattering and integrability and by Faddeev and Sklyanin [FS] where quantum inverse transformation methods were developed as a method for studying quantum integrable systems [FT]. On the basis of these developments Drinfeld [Dr1] and Jimbo [J] introduced the Hopf algebras  $U_q(\mathcal{G})$  which can be considered as deformations of universal enveloping algebras of the Kac-Moody algebras. The concept of quantization of Lie groups and Lie algebras was presented by Drinfeld in his address to Berkeley ICM [Dr1]. The construction of quantum groups based on a given solution of the Yang-Baxter equation was presented in [FRT]. The general algebraic framework of quantum groups in terms of algebras with quadratic relations [S] was developed by Manin [M].

Applications of quantum groups to low dimensional topology are my subject today. The first results in this direction were obtained in the pioneering work by Jones [Jo1] where he found a new invariant of links using certain constructions from operator algebras. Shortly after this invariant was generalized by a group of authors [HOMFLY] which is abbreviated now as HOMFLY. Another invariant similar to this was proposed by Kauffman [Ka]. The HOMFLY invariant is related to Hecke algebra, the Kauffman invariant to Birman-Wenzl algebra [BW]. Some results were generalized for invariants of graphs in  $\mathbb{R}^3$  ([Mil]).

The relation of these invariants to solutions of the Yang-Baxter equation was explained by Jones [Jo] and Turaev [Tu1]. The invariant of links related to quantum groups  $U_q(\mathcal{G})$  for any simple  $\mathcal{G}$  were studied in [Re1], where several new invariants were found. The relation between quantum groups and Hecke and Birman-Wenzl algebra was studied in [J1, Re1]. Finally it becomes clear that all these invariants can be generalized to be considered as invariants of framed graphs in  $\mathbb{R}^3$  ([RT1]).

---

\* On leave of absence, LOMI, Fontanka 27, Leningrad, 191011, USSR

Quantum groups produce interesting examples of tensor categories with non-trivial square of commutativity morphism (quasitensor categories). Algebras  $U_q(\mathcal{G})$  when  $q$  is a root of 1 give examples of semisimple tensor categories without fiber functor. The categorial explanation of invariants of links related to quantum group were given in [Tu2, TR1, Re2]. It is based on the notion of category of tangles [JS]. The general fact is that for any given tensor category [McL, DM] (with some special required properties which plays the role of ribbon element in ribbon Hopf algebra [RT1]) one can define invariants of framed links associated with these categories.

The important fact about invariants of framed links associated with quantum groups (or with tensor categories) is [RT2] that one can formulate simple conditions under which these invariants produce the invariant of framed links which is invariant under Kirby moves [Ki]. The nontrivial fact is that the algebras  $U_q(\mathcal{G})$  satisfy these conditions. It was proven for  $\mathcal{G} = sl_2$  in [RT2].

The conditions mentioned above can also be formulated as conditions on corresponding tensor category. The categorial language seems most natural for the explanation of relations between these invariants and those which arise from topological field theory. For simplest values of  $q : q^3 = 1, q^4 = 1, q^6 = 1$  invariants of 3-manifolds related to  $sl_2$ -algebra were computed in [KIM]. It was shown there that for these values of  $q$  the invariant can be reduced to known invariants. It seems that  $q^5 = 1$  should be the first value when it is not so.

Topological field theory is another important branch of theoretical physics, close to the theory of integrable systems [A] [Sh] [Wi]. The most impressive results in the conception of topological field theory were the quantum field theoretical description of Donaldson's invariants, the description of new invariants of 3-manifolds, and links in 3-manifolds obtained by Witten [Wi, Wi1].

The relation between these invariants and those mentioned above is very natural in the language of category theory and essentially based on the concept of modular category of G. Segal [Se].

## References

- [A] Atiyah, M.: Topological field theories. *Publ. Math. IHES* **68** (1989) 175–186
- [Ba] Baxter, R.J.: *Ann. Phys.* **70**, no. 1 (1972) 193–228
- [DM] Delinge, P., Milne, J.: Tannakian categories. *Lecture Notes in Mathematics*, vol. 900. Springer, Berlin Heidelberg New York 1981
- [Dr1] Drinfeld, V.G.: Quantum groups. *Proc. Intern. Congr. of Math.*, vol. 1, pp. 798–820. Amer. Math. Soc., Berkeley 1986
- [FRT] Faddeev, L.D., Reshetikhin, N.Yu., Takhtajan, L.A.: Quantization of Lie groups and Lie algebras. *Algebraic Anal.*, vol. 1, Academic Press, 1988, p. 129; *Algebra and Analysis* **1** (1989) (in Russian)
- [FS] Faddeev, L.D., Sklyanin, E.K.: *Sov. Phys. Doklady* **23** (1978) 902
- [FT] Faddeev, L.D., Takhtajan, L.A.: *Russ. Math. Surv.* **34** (1979) 11
- [HOMFLY] Freyd, P., Yetter, D., Hoste, J., Lickorish, W.B.R., Millett, K., Ocneanu, A.: A new polynomial invariant of knots and links. *Bull. Amer. Math. Soc.* **12** (1985) 239–246
- [J] Jimbo, M.: *Lett. Math. Phys.* **10** (1985) 63–70
- [J1] Jimbo, M.: *Lett. Math. Phys.* **11** (1986) 247–252
- [Jo] Jones, V.F.R.: *Bull. Amer. Math. Soc.* **12** (1985) 103–111
- [Jo1] Jones, V.F.R.: Talk at Atiyah Seminar, fall 1987 (unpublished)

- [JS] Joyal, A., Street, R.: Braided tensor categories. Preprint 1989
- [Ka] Kauffman, L.M.: An invariant of regular isotopy. *Trans. Amer. Math. Soc.* **130**, no. 2 (1990) 417–471
- [Ki] Kirby, R.: *Invent. math.* **45** (1978) 35–56
- [KiM] Kirby, R., Melvin, P.: On the invariant of Witten and Reshetikhin-Turaev for  $sl(2, \mathbb{C})$ . Preprint 1990
- [Ko] Kohno, T.: Invariants of 3-manifolds and representations of mapping class groups. Preprint 1990
- [KW] Kac, V.G., Wakimoto, M.: *Adv. Math.* **70** (1988) 156–236
- [LaS] Larson, R.G., Seedler, M.E.: *Amer. J. Math.* **141**, no. 1 (1969) 75–94
- [Lic] Lickorish, W.B.R.: *Ann. Math.* **76** (1962) 531–540
- [Lic1] Lickorish, W.B.R.: Variants of 3-manifold invariants. Preprint 1990
- [Lu] Lusztig, G.: Quantum groups at roots of 1. MIT preprint, Oct. 1989
- [M] Manin, Yu.I.: Quantum groups and non-commutative geometry. Université de Montréal preprint, CRM-1561, 1988
- [McL] MacLane, S.: Categories for the working mathematicians. Springer, Berlin Heidelberg New York 1971
- [Mi] Millett, K.C.: An invariant of 3-valent spatial graphs. UCSB preprint 1989
- [Mu] Murakami, J.: *Osaka J. Math.* **26** (1989) 1–55
- [R] Radford, D.: *Amer. J. Math.* **98**, no. 2 (1976) 333–355
- [Re1] Reshetikhin, N.Yu.: Quantized universal enveloping algebras. The Yang-Baxter equation and invariants of links. I. LOMI preprint E-4-87, 1988; II. LOMI preprint E-17-87, 1988
- [Re2] Reshetikhin, N.Yu.: Algebra and Analysis **1**, no. 2 (1989) (in Russian)
- [RT1] Reshetikhin, N.Yu., Turaev, V.G.: *Comm. Math. Phys.* **127** (1990) 1–26
- [RT2] Reshetikhin, N.Yu., Turaev, V.G.: Invariants of 3-manifolds via link polynomials and quantum groups. M.S.R.I. preprint (April, 1989); *Invent. math.* **103** (1991) 547–597
- [S] Skyanin, E.K.: *Funct. Anal. Appl.* **17** (1983) 273
- [Se] Segal, G.: The definition of conformal field theory. St. Catherine's College preprint, Oxford 1989
- [Sh] Shwartz, A.: *Lett. Math. Phys.* **2** (1978) 247–252
- [Tu1] Turaev, V.G.: *Invent. math.* **92** (1988) 527–553
- [Tu2] Turaev, V.G.: *Izv. ANSSSR* **53** (1989) 1073–1107 (in Russian)
- [Wi] Witten, E.: *Comm. Math. Phys.* **117** (1988) 353–386
- [Wi1] Witten, E.: *Comm. Math. Phys.* **121** (1989) 351–399
- [Ya] Yang, C.N.: *Phys. Rev. Lett.* **19**, no. 23 (1967) 1312–1314
- [Z] Zamolodchikov, A., Zamolodchikov, A/B: *Ann. Phys.* **120** (1979) 253–290



# Geometry of Fermionic String

*Albert Schwarz*

Department of Mathematics, University of California at Davis, Davis, CA 95616, USA

Physicists hope that the Green-Schwarz superstring theory describes all interactions existing in the Nature. However the verification of this conjecture is connected with very difficult and very interesting mathematical problems. We consider here only some problems arising in the Polyakov approach to the fermionic string. (Fermionic string is closely related with the Green-Schwarz superstring.) We explain the connection between string theory and superconformal geometry, the origin of string measure on superconformal moduli space and analytic properties of this measure, the construction of universal moduli space and the expression of string measure in terms of super  $\tau$ -function etc. The lecture is based on the papers [1–10]. The results concerning the measure on the moduli space of  $N = 2$  superconformal manifolds are new.

## Superconformal Manifolds and Strings

Let us consider a domain  $U$  in  $(1|N)$ -dimensional complex superspace  $C^{1|N}$ . One can define  $N$ -superconformal transformation of this domain as a complex analytic transformation preserving up to multiplier the 1-form

$$\alpha = dz + \sum_i \theta_i d\theta_i. \quad (1)$$

Here  $(z, \theta_1, \dots, \theta_N)$  denote complex coordinates in  $U$  ( $z$  is even,  $\theta_1, \dots, \theta_N$  are odd).  $N$ -superconformal manifold can be defined as a manifold pasted together from  $(1|N)$ -dimensional complex superdomains by means of  $N$ -superconformal transformations. The (super)space of classes of all compact  $N$ -superconformal manifolds having genus  $p$  ( $N$ -superconformal moduli space) will be denoted by  $\mathcal{M}_p^N$ . The most important cases are  $N = 0, N = 1$  and  $N = 2$ . In the case  $N = 0$  we obtain the moduli space of conformal manifolds (or moduli space of complex curves). This moduli space arises in bosonic string theory.  $N = 1$  superconformal manifolds (or simply superconformal manifolds) arise in the superstring theory. Analogously  $N = 2$  superconformal manifolds are connected with  $N = 2$  superstrings.

Let us explain the connection of string theory with moduli spaces following the Polyakov approach. The action functional of the bosonic string can be

represented in the form:

$$S(X, g) = \langle dX, dX \rangle = \int_M g^{\alpha\beta} \partial_\alpha X^\alpha \partial_\beta X^\beta dV \quad (2)$$

where  $M$  is a 2-dimensional manifold,  $X$  denotes a string field (i.e. a map of  $M$  into  $R^D$ ),  $g = (g_{\alpha\beta})$  is a riemannian metric in  $M$ ,  $dV$  denotes the corresponding volume element and  $\langle dX, dX \rangle$  denotes the scalar square of 1-form  $dX$  with respect to the metric  $g$ . The functional (2) is invariant with respect to reparametrizations and to Weyl transformations (i.e. it remains intact if we make a diffeomorphism of  $M$  or replace the metric  $g_{\alpha\beta}$  by the metric  $\varrho g_{\alpha\beta}$  where  $\varrho$  denotes a non-vanishing real function on  $M$ ). To calculate the partition function of bosonic string we must integrate  $\exp(-S)$  over all string fields and over all riemannian metrics on compact two-dimensional surfaces. The contribution of surfaces having genus  $p$  is known as the  $p$ -loop contribution to the partition function. The action functional  $S$  is quadratic with respect to string field  $X$  and therefore the calculation of the integral over  $X$  can be reduced to the calculation of  $\det \Delta_0$  where  $\Delta_0 = d^+ d$  is the Laplace operator on the scalars. Taking into account the reparametrization invariance by means of the Faddeev-Popov trick we reduce the calculation of the  $p$ -loop contribution to the partition function to the integration of

$$(\det \Delta_0)^{-D/2} \det \Delta_{gh} \quad (3)$$

over the space  $\tilde{\mathcal{M}}_p$ . Here  $\tilde{\mathcal{M}}_p$  denotes the space of orbits of the group of diffeomorphisms in the space of riemannian metrics on two-dimensional compact surface  $M$  of genus  $p$ ,  $\det \Delta_{gh}$  denotes so called ghost determinant. The action functional  $S$  is Weyl invariant; but this is not true for the expression (3) due conformal anomaly. However for  $D = 26$  (critical dimension) conformal anomaly vanishes; this permits us to represent the  $p$ -loop contribution to the partition function for critical string in the form:

$$\int_{\mathcal{M}_p} (\det \Delta_0)^{-13} \det \Delta_{gh} d\nu \quad (4)$$

where  $\mathcal{M}_p$  denotes the space obtained from  $\tilde{\mathcal{M}}_p$  by means of factorization with respect to Weyl transformations. In such a way  $\mathcal{M}_p$  consists of classes of riemannian metrics; metrics connected by diffeomorphisms and by Weyl transformations are identified. By the construction of the integrand in (4) we have chosen one of metrics in every class. The  $\det \Delta_0$ ,  $\det \Delta_{gh}$  and the volume element  $d\nu$  depend on this choice but the integrand in (4) is well defined on  $\mathcal{M}_p$ . The integrand

$$d\mu = (\det \Delta_0)^{-13} \det \Delta_{gh} d\nu \quad (5)$$

is known as the string measure on the moduli space  $\mathcal{M}_p$ . It is evident that the space  $\mathcal{M}_p$  coincides with the moduli space of conformal manifolds  $\mathcal{M}_p^{N=0}$ . The constructions above can be generalized to the supercase. In this case one has to replace the riemannian metrics on 2-dimensional surfaces by superriemannian metrics on  $(2|2)$ -dimensional supermanifolds. By definition superriemannian metrics on a  $(2|2)$ -dimensional superdomain  $U$  is an odd vector field  $e$  on  $U$

satisfying the conditions a) (anti)commutator of  $e$  and the complex conjugate field  $\bar{e}$  is a linear combination of  $e$  and  $\bar{e}$  (i.e.  $[e, \bar{e}]_+ = \alpha e + \bar{e}\bar{\alpha}$  where  $\alpha$  is an odd function), b) the vectors  $e, \bar{e}, E = [e, e]_+, \bar{E} = [\bar{e}, \bar{e}]_+$  form a basis in the tangent space. If  $e' = \exp(i\lambda)e$  where  $\lambda$  is a real function then  $e$  and  $e'$  determine the same superriemannian metric; if  $e' = \exp(i\lambda)e$ , where  $\lambda$  is an arbitrary function, one says that  $e'$  and  $e$  are connected by Weyl transformations. The action functional of a fermionic string can be written in the form

$$S(X, e) = \langle \hat{e}X, \hat{e}X \rangle \quad (6)$$

where  $X$  denotes the string field (a map of a  $(2|2)$ -dimensional supermanifold  $M$  into  $R^D$ ) and  $\hat{e}$  denotes the first order differential operator corresponding to the superriemannian metric  $e$ . This functional is invariant under reparametrizations and Weyl transformations. Slight modification of the considerations above permits us to express the  $p$ -loop contribution to the partition function of critical fermionic string ( $D = 10$ ) in terms of an integral over the space of classes of superriemannian metrics on  $(2|2)$ -dimensional supermanifolds of genus  $p$ . (We identify two metrics connected with reparametrization or Weyl transformations). One can check that this space coincides with the moduli space  $\mathcal{M}_p^{N=1}$  of superconformal manifolds. This follows from the assertion that for every superriemannian metrics on  $(2|2)$ -dimensional manifolds  $M$  one can find a covering of  $M$  by charts  $U_i$  with complex coordinates  $(z^{(i)}, \theta^{(i)})$  in such a way that in every chart superriemannian metrics takes the form  $\hat{e}^{(i)} = \Phi^{(i)}(\partial/\partial\theta^{(i)} + \theta^{(i)}\partial/\partial z^{(i)})$  (here  $\Phi^{(i)}$  denotes a non-vanishing function). We obtain a measure  $d\mu$  on the supermoduli space  $\mathcal{M}_p^{N=1}$  (one can give an expression of  $d\mu$  in terms of determinants of superLaplacians; this expression is similar to (5)).

## Superconformal Geometry

Let us return to the consideration of  $N$ -superconformal transformations of a  $(1|N)$ -dimensional complex superdomain  $U$  with coordinates  $Z = (z, \theta_1, \dots, \theta_N)$ . The operators

$$D_i = \frac{\partial}{\partial\theta_i} + \theta_i \frac{\partial}{\partial z}$$

are called covariant derivatives. Let us suppose that by the transformation  $\tilde{Z} = f(Z)$  the operators  $D_i$  transform into operators  $\tilde{D}_i$ . The transformation  $\tilde{Z} = f(Z)$  is  $N$ -superconformal if and only if

$$\tilde{D}_i = F_{ij}(Z)D_j \quad (7)$$

(the operators  $\tilde{D}_1, \dots, \tilde{D}_N$  are linear combinations of  $D_1, \dots, D_N$ ). This follows immediately from the remark that the vector fields corresponding to the operators  $D_i$  are orthogonal to the 1-form (1). It is easy to check that the matrices  $F_{ij}(Z)$  in (7) satisfy

$$F_{ij}(Z)F_{kj}(Z) = \Phi(Z)\delta_{ik}. \quad (8)$$

In other words the matrix function  $F_{ij}(Z)$  takes on values in the group  $G = O(N, C) \times \mathbf{C}^*$  where  $O(N, C)$  denotes the complex orthogonal group,  $\mathbf{C}^*$  denotes the group of non-zero complex numbers. In the case  $N = 1$  we have  $G = \mathbf{C}^*$ ; in the case  $N = 2$  the group  $G$  is disconnected and its connected part is isomorphic to  $\mathbf{C}^* \times \mathbf{C}^*$ . The  $N = 2$  superconformal transformation is called untwisted if the matrix  $F_{ij}$  in (7) belongs to the connected component of  $G$ . Recall that  $N$ -superconformal manifold is pasted together from superdomains by means of  $N$ -superconformal transformations; if in the case  $N = 2$  all these transformations are untwisted one says that  $N = 2$  superconformal manifold is untwisted. In the case  $N = 2$  it is convenient to introduce linear combinations of covariant derivatives:  $D_{\pm} = (2)^{-1/2}(D_1 \pm iD_2)$ ; the behavior of  $D_{\pm}$  by untwisted superconformal transformations is given by

$$\tilde{D}_+ = F_+ D_+, \quad \tilde{D}_- = F_- D_- . \quad (9)$$

The transformation law of the form (1) by  $N$ -superconformal transformations can be written as

$$\tilde{\alpha} = \Phi(Z)\alpha \quad (10)$$

where  $\Phi(Z) = (\det F_{ij})^{2/N}$  for  $N \geq 1$ . (For  $N = 0$  the matrix  $F_{ij}(Z)$  has no sense, but the function  $\Phi(Z)$  is well defined.) One says that a field  $\varrho$  on conformal manifold  $M$  has type  $k$  (or that  $\varrho$  is a  $k$ -differential) if the transformation law of this field by conformal transformation connecting two charts in  $M$  is given by  $\tilde{\varrho} = \Phi^{-k}\varrho$ . The field  $\varrho$  on  $N = 1$  superconformal manifold  $M$  has type  $k$  if the transformation law of  $\varrho$  by superconformal transformations is  $\tilde{\varrho} = F^{-k}\varrho$  (in the case  $N = 1$  we have only one covariant derivative  $D$  and  $\tilde{D} = FD$ , i.e. the matrix  $F$  is simply a number.) The field  $\varrho$  on untwisted  $N = 2$  superconformal manifold has type  $(k, l)$  if its transformation law is  $\tilde{\varrho} = F_+^{-k}F_-^{-l}\varrho$ . In other words the field of type  $k$  in the case  $N = 0$  is a section of a line bundle  $K^k$ , where  $K$  is a holomorphic line bundle with transition functions  $\Phi(z)$  and in the case  $N = 1$  it is a section of a line bundle  $\omega^k$  where holomorphic line bundle  $\omega$  is defined by means of transition functions  $F(Z)$ . In  $N = 2$  case the field of type  $(k, l)$  is a section of  $\omega_+^k \omega_-^l$  where  $\omega_+$  and  $\omega_-$  are defined by means of transition functions  $F_+$  and  $F_-$  correspondingly. Note that the bundles  $K$  and  $\omega$  can be interpreted also as canonical line bundles. (Recall that for an arbitrary supermanifold one can define the canonical line bundle using as transition functions (super)Jacobians of transformations connecting different charts.)

Let us denote by  $\mathcal{A}(L)$  the space of holomorphic sections of holomorphic line bundle  $L$  over a supermanifold  $M$ . In the case  $N = 0$  we will use the notation  $\mathcal{A}_k$  for  $\mathcal{A}(K^k)$  and in the case  $N = 1$  we will use the same notation for  $\mathcal{A}(\omega^k)$ . If  $M$  is a conformal manifold (a superconformal manifold) then the cotangent space at the corresponding point of  $\mathcal{M}_p^{N=0}$  (of  $\mathcal{M}_p^{N=1}$ ) can be identified with  $\mathcal{A}_2$  (with  $\Pi\mathcal{A}_3$ ). Here  $\Pi$  denotes the parity reversion. Note that  $\mathcal{A}_2$  and  $\Pi\mathcal{A}_3$  are complex linear spaces; this permits us to introduce complex structures in  $\mathcal{M}_p^{N=0}$  and  $\mathcal{M}_p^{N=1}$ . It is well known that for  $p \geq 1$  the complex dimension of  $\mathcal{M}_p^{N=0}$  is equal to  $3p - 3$ . One can prove that  $\mathcal{M}_p^{N=1}$  is a  $(3p - 3|2p - 2)$  dimensional complex supermanifold for  $p \geq 1$ . Let us show that one can construct measures

on moduli spaces using the so called Mumford form and its generalizations. First of all we will introduce the notion of a function of weight  $k$  on bases in linear space  $E$  by means of the condition  $f(\tilde{e}) = (\det P)^k f(e)$ . Here  $e = \{e_i\}$  and  $\tilde{e} = \{\tilde{e}_i\}$  are two bases in  $E$ , and  $P$  denotes the matrix connecting  $e$  and  $\tilde{e}$ ; i.e.  $\tilde{e}_i = P_i^j e_j$ . If  $E$  is a linear superspace we use the same definition, but instead of determinants one has to consider superdeterminants (Berezinians). If  $E$  is infinite-dimensional one has to consider only admissible bases (bases connected with a standard basis in  $E$  by a matrix  $P$  having well defined determinant in some sense.) The measure in linear space  $E$  can be defined as a function of weight 1 on the bases of  $E$ . We will denote the one-dimensional linear space consisting of measures in  $E$  by  $m(E)$ . In the definitions above one can take  $E$  as a complex or real linear space; correspondingly  $m(E)$  will be complex or real too. However if  $E$  is a real linear space it is natural to modify the definition of a function of bases having weight  $k$ . Namely one has to replace  $\det P$  by  $|\det P|$  in this definition. Then the measure in  $E$  can be defined as a positive function of bases having weight 1. If  $E$  is a complex linear space and  $E^{\text{real}}$  corresponding real linear space then a hermitian metric on  $m(E)$  generates positive measures in  $E^{\text{real}}$  and in  $(E^{\text{real}})^*$ . To determine a measure in a (super)manifold  $\mathcal{M}$  one has to fix the measures in all tangent spaces (or functions of weight  $-1$  on the bases of all cotangent spaces.) Let us consider a conformal manifold  $M$  and complex linear spaces

$$\Sigma_k(M) = \mathcal{A}_k(M) \dot{+} \mathcal{A}_{1-k}(M)^* \quad (11)$$

$$\lambda_k(M) = m(\Sigma_k(M)) = m(\mathcal{A}_k(M)) \otimes m(\mathcal{A}_{1-k}(M))^*. \quad (12)$$

Mumford proved that there is a canonical isomorphism

$$\lambda_k(M) \approx \lambda_1(M)^{6k^2-6k+1}, \quad (13)$$

i.e. the one-dimensional linear space  $\lambda_k(M)$  is isomorphic to the  $(6k^2 - 6k + 1)$ -th tensor power of  $\lambda_1(M)$ . (One can consider the spaces  $\lambda_k(M)$  as fibres of a holomorphic line bundle  $\lambda_k$  over  $\mathcal{M}_p^{N=0}$ ; then (13) can be interpreted as holomorphic equivalence of line bundles  $\lambda_k$  and  $\lambda_1^{6k^2-6k+1}$ .) Using (13) for  $k = 2$  we obtain an isomorphism

$$\lambda_2(M) = \lambda_1(M)^{13}; \quad (14)$$

this isomorphism is known as the Mumford form. In other words, the Mumford form can be interpreted as a function  $\mathcal{M}(e, f|M)$  of bases in  $\Sigma_2(M)$  and  $\Sigma_1(M)$ , having weight  $-1$  with respect to the basis  $e$  of  $\Sigma_2(M)$  and weight 13 with respect to the basis  $f$  of  $\Sigma_1(M)$ . There exists a natural scalar product in  $\mathcal{A}_1(M)$  (in the space of holomorphic abelian differentials). This scalar product generates a hermitian metric in  $\lambda_1(M) = m(\mathcal{A}_1(M))$ ; using this metric and the isomorphism (14) we obtain a hermitian metric in  $\lambda_2(M)$ . If a conformal manifold  $M$  has genus greater than 1,  $\lambda_2(M)$  can be identified with  $m(\mathcal{A}_2(M))$ . Remembering that  $\mathcal{A}_2(M)$  can be considered as a cotangent space to  $\mathcal{M}_p^{N=0}$  we obtain a measure in  $\mathcal{M}_p^{N=0}, p > 1$ . Using the Mumford form  $\mathcal{M}(e, f)$  we can describe this measure as follows. Let  $t_1, \dots, t_s, s = 3p - 3$ , denote a basis in the complex tangent space at

the point of  $\mathcal{M}_p^{N=0}$  corresponding to the conformal manifold  $M$  and let  $e_1, \dots, e_s$  be the dual basis in the cotangent space  $\mathcal{A}_2(M) = \Sigma_2(M)$ . Then the real measure in  $\mathcal{M}_p^{N=0}$  can be defined by the formula

$$\mu(t_1, \dots, t_q, \bar{t}_1, \dots, \bar{t}_s) = |\mathcal{M}(e, f)|^2 \quad (15)$$

where  $f$  is an arbitrary orthonormal basis in  $\mathcal{A}_1(M)$ .

A. Voronov [11] generalized the Mumford construction to the case when  $M$  is a superconformal manifold. In this case one has to modify the definitions of  $\Sigma_k(M)$  and  $\lambda_k(M)$  as follows:

$$\Sigma_k(M) = \mathcal{A}_k(M) + \Pi \mathcal{A}_{1-k}(M)^*, \quad (16)$$

$$\lambda_k(M) = m(\Sigma_k(M)) = m(\mathcal{A}_k(M)) \otimes m(\mathcal{A}_{1-k}(M)). \quad (17)$$

Voronov proved that  $\lambda_k(M)$  is canonically isomorphic to  $\lambda_1(M)^{2k-1}$ ; in particular

$$\lambda_3(M) \approx \lambda_1(M)^5. \quad (18)$$

This isomorphism is known as a super-Mumford form; it can be interpreted as a function  $\mathcal{M}(e, f)$  having weight  $-1$  with respect to the basis  $e$  in  $\Sigma_3(M)$  and weight  $5$  with respect to the basis  $f$  in  $\Sigma_1(M)$ . We will say that  $M$  is a normal superconformal manifold if  $\mathcal{A}_0(M) \equiv \mathbf{C}$  (all holomorphic functions on  $M$  are constant). Then  $\Sigma_1(M) = \mathcal{A}_1(M)$  is provided with a natural scalar product. If the genus of  $M$  is greater than  $1$ , we have  $\Sigma_3(M) = \mathcal{A}_3(M)$  and remembering the description of the cotangent space to  $\mathcal{M}_p^{N=1}$  we obtain a measure on  $\mathcal{M}_p^{N=1}$ ,  $p > 1$ . Belavin and Knizhnik [12] proved that the measure on  $\mathcal{M}_p^{N=0}$  constructed by means of the Mumford form coincides with the string measure. A corresponding result for the measure on  $\mathcal{M}_p^{N=1}$  was proved in [3].

Let us describe similar constructions for the moduli space of untwisted  $N = 2$  superconformal manifolds. First of all one has to mention that this moduli space is isomorphic to the moduli space  $\mathcal{M}_p^{1|1}$  of all complex  $(1|1)$ -dimensional supermanifolds. (For definiteness we consider compact supermanifolds of genus  $p$ .) To explain the coincidence of these two moduli spaces it is useful to note that a  $N$ -superconformal manifold can be considered as a  $(1|N)$ -dimensional complex contact supermanifold; the converse assertion is also true. (Recall that the contact structure can be specified by means of a 1-form satisfying non-degeneracy condition; the form (1) entering in the definition of  $N$ -superconformal transformations specifies contact structure in a  $(1|N)$ -dimensional complex superdomain.) For an  $(m|n)$ -dimensional contact supermanifold one can construct a  $(2m - 1|2n)$ -dimensional contact supermanifold  $PT^*X$  by means of projectivization of the cotangent bundle  $T^*X$ . If  $X$  is a  $(1|1)$ -dimensional complex supermanifold we obtain  $(1|2)$ -dimensional contact complex supermanifold, i.e.  $N = 2$  superconformal manifold. One can check that this construction gives all untwisted  $N = 2$  superconformal manifolds. The fields on the  $(1|1)$ -dimensional manifold  $X$  can be identified with chiral fields on the  $N = 2$  superconformal manifold  $PT^*X$  (i.e. with fields satisfying the condition  $D_+ \Phi = 0$ ). Let  $L$  denote a holomorphic vector

bundle over a  $(1|N)$ -dimensional supermanifold  $R$ . We define a one-dimensional complex linear space  $\lambda(L)$  by the formula

$$\lambda(L) = m(H^0(L)) \otimes m(H^1(L))^*. \quad (19)$$

Here  $H^k(L)$  are Čech cohomology groups. It is easy to check that the spaces  $\lambda_k(M)$  defined above can be considered as  $\lambda(\omega^k)$  where  $\omega$  denotes the canonical bundle over  $M$ . To define a measure on  $\mathcal{M}_p^{1|1} = \mathcal{M}_p^{N=2}$  (untwisted) we must introduce hermitian metric in  $\lambda(\omega_+^{-1}\omega_-^{-1})$  for every untwisted  $N = 2$  superconformal manifold  $R$ . (The tangent space to  $\mathcal{M}^{1|1}$  can be interpreted as  $H^1(\omega_+^{-1}\omega_-^{-1})$ .) However  $\lambda(\omega_+^{-1}\omega_-^{-1})$  is canonically isomorphic to  $\lambda(\mathcal{O}_R)$  where  $\mathcal{O}_R$  is a trivial line bundle over  $R$ . (There exists a canonical isomorphism  $\lambda(L_1) = \lambda(L_2)$  for every two holomorphic line bundles  $L_1, L_2$  over a  $(1|N)$ -dimensional manifold  $R$  where  $N \geq 2$ .) Let us consider the  $(1|1)$ -dimensional manifold  $R'$  corresponding to  $R$ . If  $x \in R'$  we denote by  $V(x)$  a  $(1|1)$ -dimensional linear space of holomorphic functions on the  $(0|1)$ -dimensional manifold  $\pi^{-1}(x)$ . (Here  $\pi$  is a natural projection of  $R$  onto  $R'$ .) The spaces  $V(x)$  can be considered as fibres of a  $(1|1)$ -dimensional vector bundle  $V$  over  $R'$ ; one can check that  $H^i(V) = H^i(\mathcal{O}_R)$  and therefore  $\lambda(V) = \lambda(\mathcal{O}_R)$ . From the other side one can identify  $\lambda(V)$  with  $\lambda(\mathcal{O}_{R'}) \otimes \lambda(\Pi\omega_{R'}) = \lambda(\mathcal{O}_{R'})^2$ . (The bundle  $\mathcal{O}_{R'}$  can be considered as a subbundle of  $V$  and the corresponding quotient bundle is  $\Pi\omega_{R'}$ .) In the normal case ( $H^0(\mathcal{O}_{R'}) = C$ ) the natural scalar product in  $\Pi H^1(\mathcal{O}_{R'}) = H^0(\omega_{R'})$  induces a hermitian metric in  $\lambda(\mathcal{O}_{R'})$  and therefore in  $\lambda(\omega_+^{-1}\omega_-^{-1}) = \lambda(\mathcal{O}_{R'})^2$ .

## Universal Moduli Space

Let us consider the space  $H$  of square integrable functions on the supercircle (i.e.  $H$  consists of functions  $F(z, \theta) = f(z) + \varphi(z)\theta$ , where  $z$  is a complex number,  $|z| = 1$  and  $\theta$  is an odd variable). We introduce an odd bilinear scalar product (bilinear pairing) in  $H$  by the formula

$$\langle F, \tilde{F} \rangle = \oint F(z, \theta) \tilde{F}(z, \theta) dz d\theta = \oint f(z) \tilde{\varphi}(z) dz + \oint \varphi(z) \tilde{f}(z) dz. \quad (20)$$

The functions  $z^n, z^n\theta$  form the standard basis of  $H$ . (Here  $n$  is an integer.) The subspace of  $H$  spanned by the vectors corresponding to  $n < 0$  (to  $n \geq 0$ ) will be denoted by  $H_-$  ( $H_+$ ). The natural projection of  $H$  onto  $H_-$  ( $H_+$ ) will be denoted by  $\pi_-$  ( $\pi_+$ ). We will say that a linear subspace  $W$  of  $H$  belongs to the super-Grassmannian  $\text{Gr}$ , if the projection  $\pi_-$  of  $W$  into  $H_-$  is a Fredholm operator and the projection  $\pi_+$  of  $W$  into  $H_+$  is a compact operator. (Recall that in the bosonic case the Grassmannian  $\text{Gr}$  can be defined in a similar way but the role of  $H$  is played by the space of functions on the circle  $|z| = 1$ .) Let us denote by  $\Gamma$  the supergroup of even invertible functions on the supercircle. The group  $\Gamma$  acts in  $\text{Gr}$  by means of multiplication operators. The subspace of  $\text{Gr}$  consisting of elements  $W \in \text{Gr}$  satisfying  $W^\perp = FW$ ,  $F \in \Gamma$  will be denoted by UMS. (Here  $W^\perp$  denotes the orthogonal complement to  $W$  with respect to the bilinear pairing (20).) All moduli spaces can be embedded in UMS by means of

the Krichever construction. More precisely, let us consider a (1|1)-dimensional complex manifold  $N$ , a point  $n \in N$  and a coordinate system  $(z, \theta)$ ,  $|z| \leq 1$ , in the neighbourhood  $U$  of  $n$ . The space of triples  $(N, n, (z, \theta))$  will be denoted by  $\mathcal{P}$ . For every point  $P \in \mathcal{P}$  we define a space  $W(P)$  as a space of functions on the supercircle admitting holomorphic extension to  $N \setminus U$  (to the exterior of supercircle). One can prove that  $W(P) \in \text{UMS}$  for every  $P \in \mathcal{P}$  and therefore  $\mathcal{P}$  is embedded in UMS. Let us consider for every element  $W \in \text{Gr}$  finite-dimensional spaces  $\mathcal{A}(W)$ ,  $\mathcal{A}(W^\perp)$  and  $\Sigma(W)$  defined by the formulas

$$\mathcal{A}(W) = W \cap H_+ = \text{Ker } \pi_-^W \quad (21)$$

$$\mathcal{A}(W^\perp) = W^\perp \cap H_+ = \Pi(H_- / \text{Im } \pi_-^W) \quad (22)$$

$$\Sigma(W) = \mathcal{A}(W) + \Pi \mathcal{A}(W^\perp). \quad (23)$$

Here  $\pi_-^W$  denotes the projection  $\pi_-$  considered as a map from  $W$  into  $H_-$  and  $\Pi$  denotes the parity reversion. If  $W = W(P, L)$ ,  $P = (N, n, (z, \theta)) \in \mathcal{P}$ ,  $L$  is a holomorphic line bundle over  $N$ , we can identify  $\mathcal{A}(W)$  with the space  $H^0(L)$  of holomorphic sections of  $L$  over  $N$ . In particular if  $N$  is a superconformal manifold,  $L = \omega^k$ , the space  $\mathcal{A}(W, P, \omega^k)$  coincides with the space  $\mathcal{A}_k(N)$  of holomorphic sections of the bundle  $\omega^k$  over  $N$  (i.e. with the space of holomorphic fields of type  $k$ ). The space  $\Sigma(W(P, \omega^k))$  can be identified with the space  $\Sigma_k(N)$ . We will construct the extension of the super-Mumford form to UMS as a function  $\mathcal{M}(w, w', W)$ , where  $w$  denotes a basis in  $\Sigma(F^3 W)$  and  $w'$  denotes a basis in  $\Sigma'(FW) = \Sigma(W^\perp) = \Sigma(W)$ . (Recall that for  $W \in \text{UMS}$  we have  $W^\perp = FW$ ,  $F \in \Gamma$ .) The construction of this extension is based on the notion of super  $\tau$ -function. Let us first define the  $\tau$ -function  $\tau(W, P)$ , where  $W \in \text{Gr}$ ,  $F \in \Gamma$  in the case when the projections  $\pi_-^W$  and  $\pi_-^{FW}$  are isomorphisms. In this case we can consider the counterimages  $(\pi_-^W)^{-1}e$  and  $(\pi_-^{FW})^{-1}e$  of the standard basis  $e = \{z^n, z^n\theta, n < 0\}$  in  $H_-$ . The  $\tau$ -function  $\tau(W, F)$  can be defined as the determinant of the matrix connecting two bases  $F(\pi_-^W)^{-1}e$  and  $(\pi_-^{FW})^{-1}e$  in  $FW$ . One can prove that this determinant is well-defined. (Note that the corresponding determinant in the bosonic case is ill-defined; therefore the definition of the super  $\tau$ -function is simpler than the definition of Sato's  $\tau$ -function.) In the general case, one can assign to every basis  $w$  of  $\Sigma(W)$  a basis  $\hat{w}$  of  $W \in \text{Gr}$ , determined up to the unimodular transformation. For example, if  $\mathcal{A}(W^\perp) = 0$ ,  $\Sigma(W) = \mathcal{A}(W)$  to construct the basis  $\hat{w}$  we add to the basis  $w$  a set  $u$  of the vectors in  $W$  satisfying the condition  $\pi_-^W(u) = e$ . (In other words the set  $u$  is transformed into the standard basis  $e$  of  $H_-$  by the projection  $\pi_-$ .) We define the  $\tau$ -function  $\tau(w, w', W, F)$  where  $w$  is the basis in  $W$ ,  $w'$  is a basis in  $FW$  as a determinant of the matrix connecting the bases  $F\hat{w}$  and  $\hat{w}'$  in  $FW$ .

$$\tau(w, \hat{w}, W, F) = \det(F\hat{w}|\hat{w}'). \quad (24)$$

This infinite-dimensional determinant is well-defined. We can now construct the extension of the super-Mumford form to UMS by the formula

$$\mathcal{M}(w, \hat{w}, W) = \frac{\tau(\Pi w', w, W, F^3)}{\tau(\Pi w', w', W, F)^3}. \quad (25)$$

Here  $W \in \text{UMS}$ ,  $F \in \Gamma$ ,  $w$  is a basis in  $\Sigma(F^3 W)$ ,  $w'$  is a basis in  $\Sigma(FW)$  and  $\Pi w'$  is a basis in  $\Sigma(W) = \Pi \Sigma(FW)$ . It is easy to check that the weights of  $\mathcal{M}(w, w', W)$  with respect to  $w$  and  $w'$  are correct ( $-1$  and  $5$  correspondingly). One can prove that for the case  $W = W(P)$ ,  $P = (N, n, (z, \theta)) \in \mathcal{P}$ ,  $N$  is a superconformal manifold the function  $\mathcal{M}(w, w', W)$  coincides with the usual super-Mumford form  $M(w, w', N)$ . Using the expression of the super  $\tau$ -function through the Reidemeister torsion we can express  $\mathcal{M}(w, w', W(P))$  through holomorphic fields on  $N$  and their zeroes. The expression obtained in such a way coincides with the similar expression for  $M(w, w', N)$  given in [11]. This remark gives the simplest proof of the relation

$$\mathcal{M}(w, w', W(P)) = M(w, w', N). \quad (26)$$

Note that the function  $\mathcal{M}(w, w', W(P))$  gives an extension of the super-Mumford form to  $\mathcal{P}$ . Moreover, it is easy to check that this function depends only on the complex manifold  $N$  (i.e. does not depend on the point  $n \in N$  and on the coordinate system  $(z, \theta)$ ). Therefore the super-Mumford form can be considered as a function on the moduli space  $\mathcal{M}^{1|1}$  of compact  $(1|1)$ -dimensional complex manifolds. (More rigorously this function depends on the point  $N \in \mathcal{M}^{1|1}$  and on the bases  $w, w'$  in  $\Sigma_3(N) = \Gamma(\omega^3)$  and in  $\Sigma_1(N) = \Gamma(\omega)$ .) The same assertion can be obtained from [11]. The moduli space  $\mathcal{M}^{1|1}$  coincides with the moduli space of untwisted  $N = 2$  superconformal manifolds and therefore the statement above permits us to assert that the super-Mumford form admits hidden  $N = 2$  superconformal symmetry; see [10]. It is well known that the  $N = 2$  world-sheet superconformal symmetry is related with the  $N = 1$  space-time symmetry in the string theory. The remarks above show that the  $N = 2$  superconformal symmetry may play a fundamental role in the string theory.

## References

1. Baranov, M., Schwarz, A.: Multiloop contribution to string theory. Pis'ma ZhETF 42 (1985) 340 (Russian). [English transl.: JETP Lett. 42 (1986) 419]
2. Baranov, M., Frolov, I., Schwarz, A.: Geometry of two-dimensional superconformal field theories. Teor. Mat. Fiz. 70 (1987) 92 (Russian). [English transl.: Teor. Math. Fiz. 70 (1987)]
3. Baranov, M., Schwarz, A.: On the multiloop contribution to the string theory. Int. J. Mod. Phys. A3 (1987) 28
4. Schwarz, A.: Fermionic string and universal moduli space. Pis'ma ZhETF 46 (1987) 340 (Russian). [English transl.: JETP Lett. 46 (1988) 428]
5. Rosly, A., Schwarz, A., Voronov, A.: Geometry of superconformal manifolds. CMP 119 (1988) 129
6. Rosly, A., Schwarz, A., Voronov, A.: Superconformal geometry and string theory. CMP 120 (1989) 437
7. Schwarz, A.: Fermionic string and universal moduli space. Nucl. Phys. B317 (1989) 323
8. Dolgikh, S., Schwarz, A.: SuperGrassmannians, super- $\tau$ -functions and strings. In: V. Knizhnik memorial volume (Brink, L., Polyakov, A. eds.). 1990
9. Dolgikh, S., Rosly, A., Schwarz, A.: Supermoduli spaces. ICTP preprint (1989), CMP 135 (1990) 91

10. Baranov, M., Frolov, I., Schwarz, A.: Geometry of the superconformal moduli space. *Teor. Mat. Fiz.* **72** (1989) 241 (Russian)
11. Voronov, A.: A formula for the Mumford measure in superstring theory. *Funk. Anal. i Prilozhen.* **22** (2) (1988) 67 (Russian). [English transl.: *Funct. Anal. Appl.* **22** (1988) 139]
12. Belavin, A., Knizhnik, V.: Algebraic geometry and the geometry of quantum strings. *Phys. Lett.* **168B** (1986) 201

# Geometric Aspects of Quantum Field Theory

Graeme Segal

Mathematical Institute, 24-29 St. Giles, Oxford OX1 3LB, England

Quantum field theory has been the basic tool of particle physics for more than half a century, but unlike earlier such tools it has not been accompanied by a satisfying mathematical theory. Recently this has begun to change. One reason is that the ideas of quantum field theory have turned out to shed light on purely mathematical questions. These applications are my subject today. So far, nevertheless, the field theory has played either a heuristic or an explanatory role in the mathematics, and the actual theorems can, and often must, be proved by other means. I hope that this will be less true as the mathematics of field theory becomes better developed. Meanwhile I shall just indicate some areas where field theory and geometry have come together, trying to illustrate the point of view rather than formulate theorems. For the most part I shall be summarizing other people's work, predominantly Witten's.

## § 1. The Framework

In his address to the Berkeley ICM Witten described  $d + 1$  dimensional quantum field theory as follows. One considers “fields” defined on some class of oriented  $d + 1$  dimensional manifolds  $M$ . A “field” might mean a map from  $M$  to some auxiliary manifold  $X$ , or a section of some natural fibre bundle on  $M$ , or even an equivalence class of such sections. In any case one has a space  $F(M)$  of fields for each compact manifold  $M$  with boundary. A field  $f \in F(M)$  has a boundary value  $f|_{\partial M}$  which belongs to some space  $F_0(\partial M)$  of fields on the boundary. We also suppose given an “action” functional  $S : F(M) \rightarrow \mathbb{R}$ , defined uniformly for all  $M$ . Then field theory is the study of the functions  $\Psi_M$  on  $F_0(\partial M)$  of the form

$$\Psi_M(f_0) = \int_{F(M; f_0)} e^{-S(f)} \mathcal{D}f, \quad (1.1)$$

where  $F(M; f_0) = \{f \in F(M) : f|_{\partial M} = f_0\}$ . More generally, if the boundary  $\partial M = \bar{\Sigma}_0 \amalg \Sigma_1$  consists of an incoming part  $\Sigma_0$  and an outgoing part  $\Sigma_1$  we are interested in operators  $\Psi_M : H_{\Sigma_0} \rightarrow H_{\Sigma_1}$ , where  $H_{\Sigma_i}$  is a space of functions on  $F_0(\Sigma_i)$ , on the form

$$(\Psi_M \phi)(f_1) = \int_{F_0(\Sigma_0)} K_M(f_0, f_1) \phi(f_0) \mathcal{D}f_0,$$

where

$$K_M(f_0, f_1) = \int_{F(M; f_0, f_1)} e^{-S(f)} \mathcal{D}f. \quad (1.2)$$

(A boundary component is ‘outgoing’ or ‘incoming’ according as its orientation agrees or not with that of  $M$ ; and  $\bar{\Sigma}_0$  denotes  $\Sigma_0$  with reversed orientation.)

The preceding formulae are only schematic, and so far it has proved impossible to develop an integration theory of the type needed. But let us at least try to abstract the essential structure. It comprises

- (i) a vector space  $H_\Sigma$  for each closed oriented  $d$ -dimensional manifold with whatever structure is appropriate;
- (ii) a bilinear pairing  $H_{\bar{\Sigma}} \times H_\Sigma \rightarrow \mathbb{C}$ ;
- (iii) an element  $\Psi_M \in H_{\partial M}$  for each  $d + 1$  dimensional manifold  $M$  with appropriate structure.

The most obvious properties these data should have are

$$(a) \quad H_{\Sigma_1 \amalg \Sigma_2} = H_{\Sigma_1} \otimes H_{\Sigma_2}$$

and

$$\Psi_{M_1 \amalg M_2} = \Psi_{M_1} \otimes \Psi_{M_2};$$

(in particular  $H_\Sigma = \mathbb{C}$  when  $\Sigma = \emptyset$ , and so  $\Psi_M \in \mathbb{C}$  when  $M$  is closed.)

(b) if two components  $\Sigma_1$  and  $\Sigma_2$  of the boundary of  $M$  are sewn together by an orientation-reversing diffeomorphism to form a new manifold  $\check{M}$  such that  $\partial M = \Sigma_1 \amalg \Sigma_2 \amalg \partial \check{M}$  then the map  $H_{\partial M} \rightarrow H_{\partial \check{M}}$  defined by the bilinear pairing takes  $\Psi_M$  to  $\Psi_{\check{M}}$ .

In particular, when  $\partial \check{M} = \emptyset$  and  $\Psi_M$  is regarded as an operator  $H_{\Sigma_1} \rightarrow H_{\Sigma_2}$ , property (b) asserts that

$$\text{trace}(\Psi_M) = \Psi_{\check{M}} \in \mathbb{C}.$$

I do not know how far an axiomatization of this kind is appropriate or helpful in traditional quantum field theory, but with some especially simple kinds of theory it works well and is a useful tool in geometry, rather like a new kind of cohomology theory. I shall mention some limitations of the framework in §§ 5 and 6 below.

A feature of each of the examples I shall describe is that either the phase space is finite dimensional because of the presence of a large group of gauge symmetries, or else the path integral (1.1) reduces to a finite dimensional integral because the integrand is an exact differential form outside a finite dimensional submanifold of the space of fields (i.e. “the stationary phase calculation is exact”). One might take this to mean that genuine quantum field theory is not involved. A more optimistic moral, however, is that one can sometimes best study finite dimensional problems by the infinite dimensional methods of field theory.

## § 2. Index Theory and the Elliptic Genus

A particle moving in a Riemannian manifold  $X$  affords the simplest example of the path integral idea, and can be regarded as a  $0 + 1$  dimensional field theory. The action for a path  $\gamma : [0, T] \rightarrow X$  is  $S(\gamma) = \frac{1}{2} \int_0^T \|\gamma'(t)\|^2 dt$ . For a point  $P$  the vector space  $H_P$  is  $L^2(X)$ , and to the 1-manifold  $[0, T]$  is associated the heat operator  $e^{-T\Delta}$  in  $H_P$ , where  $\Delta$  is the Laplacian of  $X$ . The formula (1.2) is then the usual path integral representation of the heat kernel; and if we replace  $[0, T]$  by a circle  $S_T$  of length  $T$  we have a formula for  $\text{trace}(e^{-T\Delta})$  as an integral over the loop space  $\mathcal{L}_T X = \text{Map}(S_T; X)$ . This path integral does not reduce to a finite dimensional integral.

The position is different and more relevant if we replace the action  $S : \mathcal{L}_T X \rightarrow \mathbb{R}$  with the inhomogeneous differential form  $\hat{S} = S + \omega$ , where  $\omega$  is the 2-form on  $\mathcal{L}_T X$  which to two deformations  $\xi, \eta$  of a loop  $\gamma$  assigns the number

$$\omega(\gamma; \xi, \eta) = \int_0^T \langle \xi(t), \eta'(t) \rangle dt.$$

(Here  $\eta'(t)$  is the covariant derivative.) (\*) Witten observed (see [2], [15]) that the action  $\hat{S}$  corresponds to the  $0 + 1$  dimensional field theory for which  $H_P$  is the mod 2 graded space of  $L^2$  spinor fields on  $X$ , while the operator associated to  $[0, T]$  is the spinorial heat operator  $e^{-T\hat{\Delta}}$ . The graded trace (or “supertrace”)  $\text{tr}(e^{-T\hat{\Delta}})$  is now independent of  $T$ , and is the index of the Dirac operator on  $X$ . This is a topological invariant of  $X$  called its  $\hat{A}$ -genus. On the other hand the top degree component of the differential form  $e^{-\hat{S}}$  on  $\mathcal{L}_T X$  is exact outside the finite dimensional manifold of point loops, so by Stokes’s theorem the path integral can be reduced to an integral over  $X$  (identified with the point loops). The outcome is the Atiyah-Singer formula for the index of the Dirac operator. The elaboration of this idea was described by Bismut [8] at the Berkeley ICM.

So far we have been dealing with well-known material. But we can go on to consider a  $1 + 1$  dimensional theory whose action is a differential form on the space  $F(\Sigma)$  of maps from a surface  $\Sigma$  to  $X$ . Then the vector space  $H_S$  associated to a circle  $S$  will be the space of  $L^2$  spinors on the loop space  $\mathcal{L}X$ . When  $\Sigma$  is a torus the path integral over  $F(\Sigma)$  is called the *elliptic genus*  $e_\Sigma(X)$  of  $X$ . (In fact there are a number of variants, applying in slightly different situations.) The  $0 + 1$  dimensional result that the supertrace of  $e^{-T\Delta}$  was independent of  $T$  has the analogue that  $e_\Sigma(X)$  depends only on the *conformal* structure of  $\Sigma$ , i.e. for each  $X$  it is a modular function on the upper  $\frac{1}{2}$ -plane. The elliptic genus can be interpreted formally as the equivariant index of a version of the Dirac operator on the manifold  $X$ . As with the  $\hat{A}$ -genus the path integral defining  $e_\Sigma(X)$  collapses to an integral over  $X$ , and this expresses it in terms of the characteristic numbers of  $X$  already familiar in algebraic topology. Nevertheless the elliptic genus has striking and unexpected properties, especially in connection with the topology of circle actions, and it stimulated the discovery of elliptic cohomology, a new theory whose true nature remains obscure. (An account of this subject can be found in [22]. Cf. also [1], [26]).

---

\* More accurately, we replace  $e^{-S(\gamma)}\mathcal{D}\gamma$  by  $e^{-\hat{S}}$ . The integral of an inhomogeneous form means the integral of its component of top degree.

### § 3. Topological Field Theories

A field theory is *topological* if it is defined for smooth manifolds with no additional structure (apart from a question of fixing projective multipliers which I shall suppress in this talk.) The vector spaces  $H_\Sigma$  must then be finite dimensional. A discussion of the formal properties can be found in [4].

(a) *1 + 1 dimensional Theories.* These are completely described by giving a commutative ring  $A$  with a 1 together with a linear map  $\theta: A \rightarrow \mathbb{C}$  such that the bilinear form  $(a, b) \mapsto \theta(ab)$  is non-degenerate. In fact  $A = H_{S^1}$ , and the product  $A \otimes A \rightarrow A$  is  $\Psi_M$ , where  $M$  is a disc with two holes.

(b) *2 + 1 dimensional Theories.* These are by far the most studied, and the structure is much richer: it appears to be roughly equivalent to a quantum group. A theory gives us an invariant for each closed 3-manifold, and a representation of the mapping class group of each closed surface. If we choose an element  $\xi \in H_{S^1 \times S^1}$  we get an invariant  $\phi_\xi(K, M) \in \mathbb{C}$  for each knot  $K$  in a 3-manifold  $M$  by defining

$$\phi_\xi(K, M) = \langle \xi, \Psi_{M-U} \rangle,$$

where  $U$  is a tubular neighbourhood of  $K$ . In Reshetikhin's, Turaev's, and Feigin's talks at this Congress we heard how the same output arises from a quantum group. The relation between the two approaches does not seem completely understood, but I shall say a little more about it in § 4 below.

$2 + 1$  dimensional theories are important because there is a supply of natural examples which lead to the knot invariants of Vaughan Jones and others. There is a theory for each compact Lie group  $G$  and choice of "level"  $k$ . The level is an element  $k \in H^4(BG; \mathbb{Z})$ , i.e. an integer if  $G$  is simple and simply connected. Regarding  $k$  as a characteristic class for  $G$ -bundles there corresponds to it a secondary Chern-Simons characteristic class  $S_k$  with values in  $\mathbb{R}/2\pi\mathbb{Z}$  which is defined on the space  $F(M)$  of isomorphism classes of  $G$ -bundles with connection on a 3-manifold  $M$ . This, or rather  $iS_k$ , is the action defining the theory [33]. But the theory can be constructed without mentioning path-integrals in the following way.

The vector space  $H_\Sigma$  associated to a surface  $\Sigma$  is the "quantization" of the symplectic manifold  $\mathcal{M}_\Sigma$  of flat  $G$ -bundles on  $\Sigma$ . (This is the symplectic quotient [5] of the space of *all* connections on  $\Sigma$  by the action of the gauge group.) The symplectic structure of  $\mathcal{M}_\Sigma$  depends on the level: its class is the image of  $k$  under the transgression  $H^4(BG) \rightarrow H^2(\mathcal{M}_\Sigma)$ . To obtain a definite quantization one method is to

- (i) choose a complex structure on  $\Sigma$ ,
- (ii) identify  $\mathcal{M}_\Sigma$  with the moduli space of stable holomorphic  $G$ -bundles on  $\Sigma$  by the Narasimhan-Seshadri theorem [23], thereby giving  $\mathcal{M}_\Sigma$  a Kähler structure,
- (iii) represent the symplectic form as the curvature of a holomorphic line bundle  $L$  on  $\mathcal{M}_\Sigma$ , and
- (iv) define  $H_\Sigma$  as the space of holomorphic sections of  $L$ .

One must show that  $H_\Sigma$  is essentially independent of the complex structure chosen. Even after that one must construct the vectors  $\Psi_M$  associated to 3-manifolds. No natural way of doing this is known, though in general terms one can say that if

$\Sigma = \partial M$  then the boundaries of flat  $G$ -bundles on  $M$  form a Lagrangian submanifold in  $\mathcal{M}_\Sigma$ , and this should define a vector in the quantization  $H_\Sigma$ . But the connection of the spaces  $\mathcal{M}_\Sigma$  with 3-manifolds was a great surprise, for they arise more obviously from  $1 + 1$  dimensional conformal theories, as we shall see below.

The application of this theory to the study of knots and 3-manifolds is discussed elsewhere at this Congress, so here I shall just emphasize that it has led to many new results about the geometry of the spaces  $\mathcal{M}_\Sigma$ , notably Verlinde's beautiful formula [30] for the dimension of the space  $H_\Sigma$ . By applying quantum field theory to  $\mathcal{M}_\Sigma$  in a slightly different way Witten has recently been led to conjecture a formula for the volume of  $\mathcal{M}_\Sigma$  in terms of  $\zeta_G(2g - 2)$ , where  $g$  is the genus of  $\Sigma$ , and

$$\zeta_G(s) = \Sigma(\dim V)^{-s},$$

the sum being over the irreducible representations  $V$  of  $G$ .

(c) *3 + 1 dimensional Theories.* For each compact group  $G$  there is an important  $3 + 1$  dimensional theory [3, 13, 17] which assigns to a closed 4-manifold  $W$  its Donaldson invariant, i.e. (roughly) the number of "instantons" on  $W$ . (An instanton is a solution of the self-dual Yang-Mills equations.) This theory was described in Floer's talk at this Congress. The vector space  $H_M$  for a 3-manifold  $M$  is the Floer cohomology group defined by applying infinite dimensional Morse theory to the space  $F(M)$  of isomorphism classes of  $G$ -connections on  $M$ , the Morse function being the circle-valued Chern-Simons form already mentioned. A 4-manifold  $W$  with boundary  $M$  has a relative Donaldson invariant in  $H_M$ . Unfortunately field theory, although strikingly exemplified here, has not so far helped much with the geometry, except insofar as it is a field-theoretic idea to study the instanton moduli spaces in terms of the space of all connections.

I should say a word about Floer cohomology. The infinite dimensional manifolds  $F$  which arise in field theory are usually *polarized*, in the sense that their tangent spaces are roughly decomposed into positive- and negative-energy halves. (Cf. [26] §4.) Floer's Morse function defines a decomposition of this kind, into the positive and negative eigenspaces of the Hessian. For such a manifold  $F$  one expects to be able to define "middle dimensional cohomology", by considering infinite dimensional cycles whose tangent spaces roughly fill the negative half of the tangent spaces to  $F$ . This idea goes back, of course, to Dirac's treatment of electrodynamics in terms of a sea of negative energy electrons. The same idea has been formalized by Feigin in his "semi-infinite" cohomology of Lie algebras [14]. Apart from the space of connections above, Floer cohomology has also been applied to the loop space of a symplectic manifold [12], and there too it arises as the state space of a field theory, the  $1 + 1$  dimensional topological  $\sigma$ -model of [32].

## § 4. 1 + 1 Dimensional Conformal Field Theory

Conformal field theory is akin to topological field theory in the sense that, up to isomorphism, a compact surface has only a finite dimensional space of conformal structures. Conformal theories can be axiomatized in the same way as topological ones [27, 28]. They have been much studied since the influential paper [7], partly

for their relevance to string theory, but also because of their role in at least three areas of mathematics:

(i) the representation theory of loop groups and of  $\text{Diff}(S^1)$ ,

(ii) the study of the moduli spaces of Riemann surfaces and holomorphic bundles, and

(iii) the construction of the monster simple group and its representations.

For the third of these areas I refer to [15]. A slightly more conventional approach to conformal field theory is summarized in [18].

A conformal theory consists of a vector space  $H$  naturally associated to the standard circle  $S^1$ , together with an operator  $\Psi_\Sigma : H^{\otimes m} \rightarrow H^{\otimes n}$  for each Riemann surface  $\Sigma$  with  $m$  incoming and  $n$  outgoing parametrized boundary circles. Thus  $\text{Diff}(S^1)$  acts (projectively) on  $H$ , and so does the semigroup  $\mathcal{A}$  of surfaces which are topologically cylinders. (The composition-law is sewing end-to-end.) The semigroup  $\mathcal{A}$  has twice the dimension of  $\text{Diff}(S^1)$ , and is a complex manifold. One of the important ideas of the theory is that  $\mathcal{A}$  plays the role of a complexification of the group  $\text{Diff}(S^1)$ : more precisely, the relation between them is the same as that between the unitary group  $U_n$  and the semigroup  $\{g \in GL_n(\mathbb{C}) : \|g\| < 1\}$  of contraction operators. (Cf. [28] and also Neretin [24].)

Let us recall that the loop group  $\mathcal{L}G$  of a compact group  $G$  has an interesting class of irreducible projective representations  $\{H_{k,V}\}$  – the positive energy representations – which are parametrized by their level  $k \in H^4(BG; \mathbb{Z})$ , which describes the projective multiplier of the representation, and an irreducible representation  $V$  of  $G$ . (The image of  $k$  in  $H^2(\mathcal{L}G)$  is the class of the circle bundle defined by the central extension.) For a given level  $k$  only a finite set of representations  $V$  can occur. An important fact about the representations  $H_{k,V}$  is that they possess a canonical projective action of  $\text{Diff}(S^1)$  intertwining with that of  $\mathcal{L}G$ . This action extends to – more accurately, is the boundary value of – an action of  $\mathcal{A}$  by trace-class operators  $\Psi_A : H \rightarrow H$ . In fact  $\Psi_A$  is characterized by intertwining with the group  $G_A$  of holomorphic maps  $A \rightarrow G_{\mathbb{C}}$ , which acts on the source and target of  $\Psi_A$  via restriction to the two ends of  $A$ . The remarkable fact is that the irreducible representations of a given level constitute something very close to a conformal field theory. To state this precisely one needs the concept of a modular functor [27, 28]. (In the literature modular functors are usually referred to as “conformal blocks” [7] or “solutions of the Knizhnik-Zamolodchikov equations”. For the latter, see Varchenko’s talk at this Congress.)

A *modular functor* has a finite set  $\Phi$  of labels. It assigns a finite dimensional vector space  $E_\Sigma$  to each Riemann surface with boundary where each boundary circle is labelled with an element of  $\Phi$ . The axioms are

(i)  $E_\Sigma = \mathbb{C}$  when  $\Sigma$  is the Riemann sphere,

(ii)  $E_{\Sigma_1 \sqcup \Sigma_2} \cong E_{\Sigma_1} \otimes E_{\Sigma_2}$ ,

(iii)  $E_{\check{\Sigma}} \cong \bigoplus_{\phi \in \Phi} E_{\Sigma, \phi}$  where  $(\Sigma, \phi)$  is obtained from  $\check{\Sigma}$  by cutting it along a simple closed curve and giving both new boundary circles the label  $\phi$ .

For the application  $\Phi$  is the set of irreducible representations of  $\mathcal{L}G$  of a given level. There is a modular functor  $E$  such that when  $\Sigma$  is a Riemann surface with  $m$

incoming and  $n$  outgoing boundary circles labelled with representations  $H_{\alpha_1}, \dots, H_{\alpha_m}$  and  $H_{\beta_1}, \dots, H_{\beta_n}$  there is an operator

$$\Psi_{\Sigma, \xi}: H_{\alpha_1} \otimes \cdots \otimes H_{\alpha_m} \rightarrow H_{\beta_1} \otimes \cdots \otimes H_{\beta_n}$$

for each  $\xi \in E_\Sigma$  which intertwines with the action of the group  $G_\Sigma$  of holomorphic maps  $\Sigma \rightarrow G_{\mathbb{C}}$ . (In fact  $E_\Sigma$  can be *defined* as the space of such intertwining operators: then the point to establish is property (iii) above, which amounts to a version of the Peter-Weyl theorem for loop groups.) The first complete proof of this result is in [29]. (Cf. Tsuchiya's talk at this Congress.)

One of the advantages of the field-theoretic viewpoint in the representation theory of loop groups is to make plain the otherwise mysterious modularity properties of the characters: in field theory the values of the characters are naturally associated to complex tori.

Witten realized [33] that the modular functor  $E_\Sigma$  just described is essentially independent of the complex structure of  $\Sigma$ , and is the state space of the corresponding 2 + 1 dimensional topological theory based on the Chern-Simons action. More recently Kontsevich [21] has sketched an argument to show, still more surprisingly, that the concepts of modular functor and 2 + 1 dimensional topological theory are exactly equivalent.

A “topological” modular functor is closely related to a quantum group, for the quantum deformation of  $G$  amounts essentially to a way of defining an exotic tensor product on the category of representations of  $G$ . For a modular functor  $E$  we can define

$$V_1 \otimes_E V_2 = \bigoplus_W E_{\Sigma, V_1, V_2, W} \otimes W,$$

where  $W$  runs through the irreducible representations of  $G$ , and  $\Sigma$  is a disc with two holes whose boundary components are labelled  $V_1, V_2$  (incoming) and  $W$  (outgoing).

It is easy to relate the modular functor  $E_\Sigma$  to the space  $H_\Sigma = \Gamma(\mathcal{M}_\Sigma; L)$  of holomorphic sections described in § 3. Let us decompose the closed surface  $\Sigma$  as  $\Sigma_1 \cup \Sigma_2$  by a simple closed curve  $S$ . A holomorphic bundle on  $\Sigma$  is automatically trivial on  $\Sigma_1$  and  $\Sigma_2$ , so it can be described by a clutching function on  $S$ , i.e. by an element of  $\mathcal{L}G_{\mathbb{C}}$ . The set of isomorphism classes of bundles on  $\Sigma$  – essentially the same as  $\mathcal{M}_\Sigma$  – is therefore the double coset space  $G_{\Sigma_1} \backslash \mathcal{L}G_{\mathbb{C}} / G_{\Sigma_2}$ , and the space  $\Gamma(\mathcal{M}_\Sigma; L)$  is the  $G_{\Sigma_1}$ -invariant part of  $H = \Gamma(\mathcal{L}G_{\mathbb{C}} / G_{\Sigma_2}; \pi^* L)$ , where  $\pi: \mathcal{L}G_{\mathbb{C}} / G_{\Sigma_2} \rightarrow \mathcal{M}_\Sigma$ . If we now take  $\Sigma_2$  to be a standard disc then  $H$  is the basic representation of  $\mathcal{L}G$  of level  $k$ , constructed by the Borel-Weil method [25]. Finally, it is easy to see that when the boundary of  $\Sigma_1$  is labelled with  $H$  we have  $E_\Sigma = E_{\Sigma_1}$ , and so

$$E_\Sigma \cong H^{G_{\Sigma_1}} \cong \Gamma(\mathcal{M}_\Sigma; L).$$

The preceding argument, which shows how representations of  $\mathcal{L}G$  define functions on the moduli space of  $G$ -bundles, also shows how representations of  $\text{Diff}(S^1)$  give functions on the moduli space  $\mathcal{C}_\Sigma$  of complex structures on a smooth surface  $\Sigma$ . For  $\mathcal{C}_\Sigma$  behaves like a double coset space of the semigroup  $\mathcal{A}$ : if we write  $\Sigma = \Sigma_1 \cup \Sigma_2$ , and choose fixed complex structures on  $\Sigma_1$  and  $\Sigma_2$ , then  $\Sigma_1 \cup A \cup \Sigma_2$  runs through an open set of  $\mathcal{C}_\Sigma$  as  $A$  runs through  $\mathcal{A}$ . The representation theory of

$\text{Diff}(S^1)$  allows us, for example, to identify and classify holomorphic line bundles on  $\mathcal{C}_\Sigma$  much more simply than does conventional algebraic geometry. (Cf. [6, 28].)

## § 5. Zamolodchikov's $c$ -Theorem

The point of view of this talk is successful with topological and conformal field theories, but so far it has never been taken seriously in a wider context. At present the only general definition of a field theory is the classical one in terms of the vacuum expectation values of a class of operators varying from theory to theory. This is not sufficiently manageable for one to be able to speak, for instance, of the “space  $\mathcal{T}$  of all  $1 + 1$  dimensional theories”. Nevertheless one of the most interesting recent developments has been the following result of Zamolodchikov [36], which is framed in terms of the space  $\mathcal{T}$ .

Whatever may be the definition of a theory, it will presumably be true that from any theory  $T$  one can derive a 1-parameter family of theories  $T_\lambda$  (for  $\lambda \in \mathbb{R}^\times$ ) simply by multiplying all lengths by  $\lambda$ . The resulting flow on  $\mathcal{T}$  is the *renormalization group flow*. Conformal theories are fixed points of this flow. Zamolodchikov's idea is to define a Riemannian metric on  $\mathcal{T}$  and a smooth function  $c : \mathcal{T} \rightarrow \mathbb{R}$  such that

- (i) the renormalization group flow is the gradient flow of  $c$ , and
- (ii)  $c(T)$  is equal to the central charge if  $T$  is a conformal theory.

The *central charge* of a conformal theory is the number describing the central extension of  $\text{Diff}(S^1)$  which acts on the state space of the theory.

It is fairly straightforward to calculate the possible non-conformal infinitesimal deformations of a conformal theory, so Zamolodchikov's theorem suggests that one could in principle discover the global topology of the space  $\mathcal{T}$  by Morse theory. Vafa and others have made some steps in this direction.

Zamolodchikov's argument is based on perturbation theory. I think it is a fascinating challenge to put it in a better mathematical setting.

## § 6. $1 + 1$ Dimensional Quantum Gravity

The most dramatic recent development in quantum field theory has been a breakthrough in 2-dimensional quantum gravity. It has suddenly appeared possible to perform integrals over the space of all metrics on a surface and get very explicit answers. The main success has come from the technique of random matrices [9, 11, 19], and I cannot discuss it here. One very unexpected outcome is to link the theory with the classical completely integrable systems of non-linear partial differential equations such as the KdV equation. So far the situation has not really been assimilated mathematically, but the results seem to describe the algebraic topology of the space of metrics on a surface, or – equivalently – the moduli spaces of complex structures.

From ordinary algebraic geometry one knows (see Morita's talk at this Congress) a ring of stable cohomology classes on the moduli spaces  $\mathcal{M}_g$  of surfaces of genus  $g$  (“stable” means that they are defined independently of  $g$ ), and Witten has

claimed that the field-theoretic results are the integrals of these classes over the spaces  $\mathcal{M}_g$ , i.e. the *characteristic numbers* of  $\mathcal{M}_g$ . This has led him to the striking conjecture [34] that the generating function for the characteristic numbers is a certain specific solution of the KdV hierarchy. His lectures [35] give an excellent account of the present state of the subject.

From the point of view of this talk it is interesting that quantum gravity does not fit directly into the framework of § 2 above, but nevertheless seems likely to be axiomatizable along related but more subtle lines. The crucial point that distinguishes the gauge-theory situations to which § 2 applies from the gravitational ones to which it does not is simply that an automorphism of a bundle  $P$  on  $M_1 \cup M_2$  is just a pair of automorphisms of  $P|M_1$  and  $P|M_2$ , but a diffeomorphism of  $M_1 \cup M_2$  cannot be broken into two diffeomorphisms.

## References

1. O. Alvarez, T. Killingback, M. Mangano, P. Windey: String theory and loop space index theorems. *Comm. Math. Phys.* **111** (1987) 1–10
2. M.F. Atiyah: Circular symmetry and stationary phase approximation. *Proc. Conf. in honour of L. Schwartz, Astérisque* **131** (1985) 43–59
3. M.F. Atiyah: New invariants of 3- and 4-dimensional manifolds. *Amer. Math. Soc., Proc. Symp. Pure Math.* **48** (1988) 285–99
4. M.F. Atiyah: Topological quantum field theories, *Publ. Math. I.H.E.S. Paris* **68** (1989) 175–86
5. S. Axelrod, S. Della Pietra, E. Witten: Geometric quantization of Chern-Simons gauge theory. *J. Diff. Geom.* **33** (1991) 787–902
6. A.A. Beilinson, V.V. Schechtman: Determinant line bundles and Virasoro algebras. *Comm. Math. Phys.* **118** (1990) 651–701
7. A.A. Belavin, A.M. Polyakov, A.B. Zamolodchikov: Infinite conformal symmetry in two-dimensional quantum field theory. *Nucl. Phys.* **B241** (1984) 333–380
8. J.-M. Bismut: Index theorem and the heat equation, *Proc. Internat. Cong. Math. Berkeley 1986*, pp. 491–504
9. E. Brezin, V.A. Kazakov: Exactly solvable field theories of closed strings. *Phys. Lett.* **236B** (1990) 144–150
10. S.K. Donaldson, P. Kronheimer: *Geometry of 4-manifolds*. Oxford U.P. 1990
11. M. Douglas, S. Shenker: Strings in less than one dimension. *Nucl. Phys. B* **335** (1990) 635–654
12. A Floer: Morse theory for fixed points of symplectic diffeomorphisms. *Bull. Amer. Math. Soc.* **16** (1987) 279–281
13. A. Floer: An instanton invariant for 3-manifolds. *Comm. Math. Phys.* **118** (1988) 215–240
14. I.B. Frenkel, H. Garland, G. Zuckerman: Semi-infinite cohomology and string theory, *Proc. Nat. Acad. Sci. USA* **83** (1986) 8442–6
15. I.B. Frenkel, J. Lepowsky, A. Meurman: *Vertex operator algebras and the monster*. Academic Press, 1988
16. D. Friedan, P. Windey: Supersymmetric derivation of the Atiyah-Singer index theorem and the chiral anomaly. *Nucl. Phys.* **B235** (1984) 395–416.
17. M. Furuta, D. Kotschick, S. Donaldson: Floer homology groups in Yang-Mills theory. (To appear)
18. K. Gawedzki: Conformal field theory. *Seminaire Bourbaki* **704** (1988) in *Astérisque*.

19. D.J. Gross, A.A. Migdal: A nonperturbative treatment of two dimensional quantum gravity. *Phys. Rev. Lett.* **64** (1990) 127–130
20. N.J. Hitchin: Flat connections and geometric quantization. *Comm. Math. Phys.* **131** (1990) 347–380
21. M. Kontsevich: Rational conformal field theory and invariants of 3-dimensional manifolds (Marseilles preprint, to appear)
22. P.S. Landweber (ed.): Elliptic curves and modular forms in algebraic topology, Proceedings, Princeton 1986. (Lecture Notes in Mathematics, vol. 1326.) Springer, Berlin Heidelberg New York 1986
23. M.S. Narasimhan, C.S. Seshadri: Stable and unitary bundles on a compact Riemann surface. *Ann. Math.* **82** (1965) 540–67
24. Yu. A. Neretin: A complex semigroup containing the group of diffeomorphisms of a circle. *Funkt. Anal. Prilozh.* **21** (1987) 82–83
25. A. Pressley, G. Segal: Loop groups. Oxford, U.P. 1986
26. G. Segal: Elliptic cohomology, *Seminaire Bourbaki* **695** (1988). In *Astérisque* **161–162** (1988) 187–201
27. G. Segal: Two dimensional conformal field theories and modular functors. IXth Proc. Internat. Cong. Math. Phys. Swansea 1988 (B. Simon, A. Truman, I.M. Davies, (eds.)) Adam Hilger, 1989, pp. 22–37
28. G. Segal: The definition of conformal field theory. (To appear)
29. A. Tsuchiya, K. Ueno, Y. Yamada: Conformal field theory on the universal family of stable curves with gauge symmetry. In *Conformal field theory and solvable lattice models*. *Adv. Stud. Pure Math.* **16** (1988) 297–372
30. E. Verlinde: Fusion rules and modular transformations in two dimensional conformal field theory. *Nucl. Phys. B* **300** (1988) 360–376
31. E. Witten: Topological quantum field theory. *Comm. Math. Phys.* **117** (1988) 353–86
32. E. Witten: Topological  $\sigma$ -models. *Comm. Math. Phys.* **118** (1988) 411–449
33. E. Witten: Quantum field theory and the Jones polynomial. *Comm. Math. Phys.* **121** (1989) 351–99
34. E. Witten: On the structure of the topological phase of two dimensional gravity. IAS Preprint IAS SNS HEP 89/66
35. E. Witten: Two-dimensional gravity and intersection theory on moduli space. Lectures, Harvard 1990 (To appear)
36. A.B. Zamolodchikov: Renormalization group and perturbation theory near fixed points in two-dimensional field theory. *Yad. Fiz.* **46** (1987) 1819–1831

# Quantum Mechanics of Many-Particle Systems\*

I.M. Sigal<sup>1</sup>

Department of Mathematics, 100 St. George Street, University of Toronto  
Toronto, Ontario, M5S 1A1

## 1. Introduction

In this talk I will discuss some mathematical questions arising in the theory of quantum many-body systems. Examples of such systems are atoms, molecules, nuclei, solids and, to some extent, stars. The remarkable fact is that such diverse objects are described within a single mathematical framework given in terms of the Schrödinger operator:

$$H = -\Delta + V(x) \quad \text{on } L^2(X).$$

Here  $X$  is the configuration space of a system in question, it is either  $\mathbb{R}^{3N}$  or a linear subspace thereof, where  $N$  is the number of particles,  $\Delta$  is the Laplacian on  $X$  and  $V(x)$  is the total potential energy, e.g.

$$V(x) = \sum V_{ij}(x_i - x_j),$$

the sum of pair potentials, with  $x_i$  being the coordinate of the  $i$ th particle and  $x = (x_1, \dots, x_N) \in X$ . Whenever the Pauli principle is taken into account,  $L^2(X)$  should be replaced by its appropriate subspace (see the paragraph after Equation (12)).

Physical properties of quantum systems are associated with spectral characteristics of  $H$ . For instance, existence and uniqueness of unitary dynamics, i.e. solution of the Cauchy problem

$$i \frac{\partial \psi_t}{\partial t} = H \psi_t \quad \text{and} \quad \psi_0 = \psi, \tag{1}$$

satisfying  $\|\psi_t\| = \|\psi\|$ , is equivalent to the statement that  $H$  is self-adjoint. Once this is established, the next problem is classification of orbits  $\psi_t$ . Though  $\psi_t$  is an orbit in the Hilbert space  $L^2(X)$ , of fundamental interest is its behaviour in the configuration space  $X$ . According to the latter, orbits are classified as bounded, wandering and escaping ( $\rightarrow \infty$ ). This rough decomposition turns out to be equivalent to splitting the spectrum into pure point, singular continuous and absolutely continuous parts. Strange, recurrent orbits coming from the singular continuous spectrum subspace are ruled out by

---

\* Supported by NSERC under Grant NA7901.

<sup>1</sup> I.K. Killam Research Fellow.

**Theorem** [Mo 1981, PSS 1981]. *Assume that*

$$|y|^{\alpha} \partial^{\alpha} V_{ij}(y) \text{ are } \Delta_y\text{-compact for } |\alpha| \leq 2. \quad (2)$$

*Then*

$$\text{sing. cont. spec.} = \phi.$$

The next basic questions are the asymptotic behaviour of unbounded orbits and stability of bounded ones under perturbations of  $H$ . They lead to the three main problems of the theory of quantum many-particle systems: scattering, binding and resonances. We consider here the first two problems, referring the interested reader to [Sig 1989] for a review of the last one.

## 2. Scattering Theory

The basic problem here is to show that as  $|t| \rightarrow \infty$ , every orbit starting in the continuous spectrum subspace obeys

$$\|\psi_t - \sum \text{simple orbits}\| \rightarrow 0. \quad (3)$$

Here the sum is taken over all possible break-ups of the system into subsystems and over all stable states (i.e. eigenfunctions of the corresponding Schrödinger operators) of the resulting subsystems. Given a break-up and a specification of stable states, the corresponding simple motion can be written as

$$\text{simple motion} = \phi \otimes u_t,$$

where  $\phi$  is the product of the eigenfunctions of the subsystems (depending on the internal coordinates) and  $u_t$ , a free (no potentials) orbit of the centers-of-mass of the subsystems (see Fig. 1). The time factors corresponding to the eigenvalues of the subsystems (the internal energies) are included into  $u_t$ .

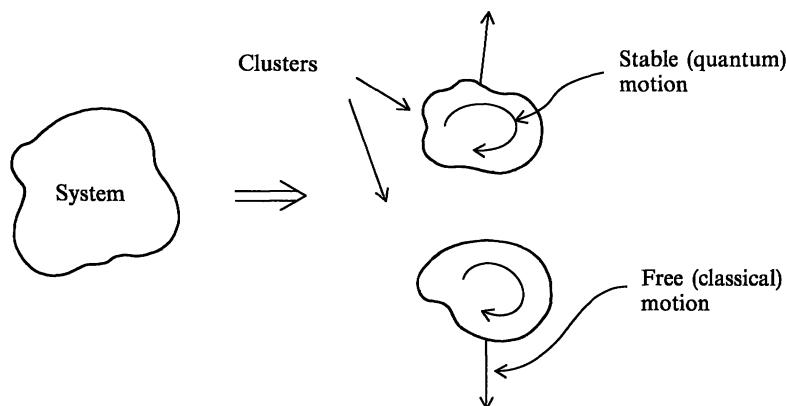


Fig. 1. Break up of a system into independent stable clusters

This problem is called *asymptotic completeness*. Over the last 40 years it was a focus of attention of many mathematicians and mathematical physicists. Important contributions were made by T. Kato, S. Kuroda, T. Ikebe, M. S. Birman, L. D. Faddeev, V. Enss, B. Simon, S. Agmon and L. Hörmander among others (see references in [SigSof 1987, 1990ab]). The problem was solved for the short-range potentials in [SigSof 1987] (see [Graf 1990] for a beautiful and concise proof, different proofs are given in [Kit 1990, Tam 1990]):

**Theorem** [SigSof 1987]. *Assume (2) and that*

$$\langle y \rangle^\mu V_{ij}(y) \quad \text{are } \Delta_y\text{-bounded} \quad (4)$$

*with  $\mu > 1$ . Then asymptotic completeness holds.*

Precise restrictions on smoothness of  $V_{ij}$  are not essential and can be relaxed if some extra care is exercised. What is important is their decay at infinity. If  $\mu > 1$ , then the potentials are called *short-range*, and if  $\mu \leq 1$ , *long-range*. Choices of the free evolutions are different for short-range and long-range cases. In particular, in the short-range case they are generated by the Hamiltonians

$$H^{\text{asympt}} = E^{\text{intern}} - \Delta^{\text{CM}}, \quad (5)$$

where  $E^{\text{intern}}$  is the sum of the eigenvalues of the subsystems and  $\Delta^{\text{CM}}$  is the Laplacian in the center-of-mass coordinates, and in the long-range case, by more complicated, time-dependent operators. For the long-range potentials only partial results are, presently, obtained. Namely, for  $N \leq 3$  and  $\mu > \sqrt{3} - 1$  asymptotic completeness is proven in [Enss 1985] and for  $N \leq 4$  and  $\mu = 1$ , in [SigSof 1990a,b].

**Open Problem.** Prove asymptotic completeness for all  $N$  and long-range (especially, Coulomb-type:  $\mu = 1$ ) potentials.

We make a few comments about the proofs.

**Propagation Set.** The first step one makes in the modern scattering theory is changing the viewpoint. Instead of studying the evolution of the system in the physical space  $\mathbb{R}^3$  one investigates the behaviour of orbits  $\psi_t$  in the phase-space  $T^*X$ . For  $\Omega \subset T^*X$  we define (modulo boundary terms) the probability that the system in question is in  $\Omega$  at time  $t$ :

$$(\text{Prob. syst. in } \Omega \text{ at } t) \equiv \|\phi\psi_t\|^2, \quad (6)$$

where  $\phi$  is a pseudodifferential operator whose symbol is supported in  $\Omega$  and is equal to 1 in a slightly smaller set. If  $\Omega \subset X$  and  $\phi$  is the multiplication operator by the characteristic function of  $\Omega$ , then (6) is a standard definition for a quantum probability for the system to be in the region  $\Omega$  of the configuration space. One expects that as  $|t| \rightarrow \infty$ ,  $\psi_t$  concentrates on the set

$$\text{PS} = \cup (\text{Class. phase-space traj. of quant. stable subsystems}),$$

called the *propagation set* (see Fig. 1). Indeed, we have

**Theorem** (SigSof 1987). *Assume (2) and (4) with  $\mu > 0$ . Then as  $|t| \rightarrow \infty$*

$$(\text{Prob. syst. in } T^*X \setminus PS \text{ at } t) \rightarrow 0 \quad (7)$$

*in the mean sense explained below.*

Note that unlike the wavefront set in the propagation of singularities the propagation set consists of *partially quantized* bicharacteristics (classical trajectories). A peculiarity of quantum propagation is that the motion along these bicharacteristics is unstable: the system can tunnel from one bicharacteristic to another. One of the consequences of this is that the convergence in (7) is in the mean sense:

$$\int_{|t| \geq 1} (\text{l.h.s. of (7)}) \frac{dt}{|t|} < \infty.$$

**Thresholds.** A key restriction, omitted from the statement of the above theorem, says that initial conditions for the orbits  $\psi_t$  must be from Range  $E_A(H)$ , where  $E_A(H)$  is the spectral projection of  $H$  and  $A$  is some interval away from the thresholds. The *thresholds* or threshold energies are critical values of Hamiltonians driving partially quantized bicharacteristics, (i.e.  $E^{\text{intern}}$  of (5)). Understanding the evolution at threshold energies is a delicate problem, particularly, because of the tunneling between critical and non-critical bicharacteristics (i.e. bicharacteristics for which a given energy is a critical value and those for which it is not).

One of the principal differences between treatments of short-range and long-range scattering is that while in the short-range case one can avoid considering threshold energies, in the long-range case, one cannot.

**Long-Range Scattering.** In the long-range case as a system breaks apart, the interaction between departing subsystems cannot be neglected entirely. It is replaced by a potential depending only on the internal coordinates of the subsystems and time. As a result the induction in the number of particles, which is used in the proof, forces one to estimate orbits  $\psi_t$  generated by time-dependent Schrödinger operators of the form

$$H(t) = H + W(x, t),$$

where the potential  $W$  obeys

$$|\partial_{(x,t)}^\alpha W(x, t)| \leq C_\alpha (1 + |x| + |t|)^{-\mu - |\alpha|}.$$

Since the energy is not conserved anymore, one cannot separate and remove threshold energies by cut-off functions in  $H$ , as it is done in the short-range case. To study  $\psi_t$  we use a version of microlocal analysis with several fine time scales. This allows us to separate critical bicharacteristics from non-critical ones. More precisely, we split the energy axis as

$$\text{dist}(E, \text{thresholds}) \geq t^{-\beta} \quad (8)$$

$$\text{dist}(E, \text{threshold}) \leq t^{-\beta} \quad (9)$$

with some  $\beta < \mu$  ( $= 1$ ). For  $E_{A(t)}(H)\psi_t$  with  $A(t)$  from (8) we still manage to prove that subsystems separate (though with relative velocities vanishing in time). For energies in (9), we consider separately regions

$$|x| \leq t^\alpha \quad (10)$$

and

$$|x| \geq \frac{1}{2}t^\alpha \quad (11)$$

with  $\alpha$  depending on the directions in  $X$  and obeying  $(1 \Rightarrow) \mu > \alpha > 1 - \beta/2 > 0$ . In the first of these regions we show that the system in question is localized. To this end we use that non-critical bicharacteristics lie outside (10), while the critical ones stay inside. Indeed, their velocities, due to (9), are bounded by

$$\begin{aligned} \text{critic. vel.} &\leq C(\text{dist}(E, \text{thresholds}))^{1/2} \\ &= Ct^{-\beta/2}. \end{aligned}$$

As a result we can set  $x = 0$  in  $W(x, t)$  which reduces the problem to the time-independent one. Since the critical bicharacteristics do not reach the second region, the system propagates there along non-critical bicharacteristics. This means it breaks up as  $|t| \rightarrow \infty$ . Subtlety here is that we can conduct the analysis above only in fixed directions, not globally, reducing the problem to a simple one step by step. Thence stems our limitation on the number of particles.

### 3. Binding

Binding is the property of matter to form stable compounds, such as atoms, molecules, nuclei, etc. It is measured by a gap between  $\inf \text{spec } H$  and  $\text{cont spec } H$ . This gap is called the *binding energy*,

$$\text{BE} = \inf \text{cont spec } H - \inf \text{spec } H.$$

In fact, this is the energy needed to destroy the most stable bounded orbit, the one corresponding to  $\inf \text{spec } H$ .

**Example: Atoms.** Consider a system consisting of  $N$  electrons and a nucleus of charge  $z$ . Let  $z = N$ . For simplicity we assume the nucleus to be infinitely heavy. Then the Schrödinger operator of such a system is

$$H = \sum_{i=1}^N (-\Delta_i + V(x_i)) + \frac{1}{2} \sum_{i \neq j} |x_i - x_j|^{-1}, \quad (12)$$

where  $V(x) = -z/|x|$  and the units used are such that the electron mass is  $1/2$ , the electron charge is  $1$  and the Planck constant is  $2\pi$ . It acts on  $\bigwedge_{i=1}^N L^2(\mathbb{R}^3 \times \mathbb{Z}_q)$ , where  $q = 2$ , the number of spin states of electron. It was shown in [LiebSim 1977] that

$$-\inf \text{spec} = O(z^{7/3}).$$

On the other hand it is believed that

$$\text{BE} = O(1).$$

Thus the existence of matter as we know it is due to a subtle phenomenon, indeed.

**Theorem** [SecoSigSol 1990].

$$\text{BE} \leq \text{const } z^{20/21}. \quad (13)$$

The proof of this theorem is rather instructive. Some of the tools used are described in the following sections.

**Open Problem.** Show that  $\text{BE} = O(1)$ .

Another way to measure binding is by the maximal number of electrons a nucleus of charge  $z$  can hold together,  $N(z)$ . It was shown in [Rus 1982, Sig 1982] that  $N(z)$  is finite.

**Theorem.**  $N(z) = z + O(z^{5/7})$ .

It was shown in [BengLieb 1983, Solov 1990] that if the electrons were bosons, then

$$N^{\text{boson}}(z) = 1.21z + O(z)$$

(the numerical value of the coefficient in front of  $z$  was computed in [Baum 1984]). [Bach 1990] proves that

$$\text{BE}^{\text{boson}} = O(z^2).$$

Thus Pauli principle plays a crucial role in the formation of atoms.

The asymptotic behaviour of  $N(z)$  was established in [LSST 1988] and the remainder estimate was derived in [FeffSeco 1990a]. A simpler proof is given in [SecoSigSol 1990].

**Ground State Energy.** One of the basic characteristics of a quantum system is its ground state energy, i.e. the lowest eigenvalue of the corresponding Schrödinger operator. As an example we consider a molecule with  $N$  electrons and  $M$  nuclei of charges  $z_1, \dots, z_M$ . Let  $\sum z_i = N$ . We assume the nuclei to be infinitely heavy and located at positions  $r_1, \dots, r_M$  (Born-Oppenheimer model). Let  $Z = (z_1, \dots, z_M)$ ,  $|Z| = \sum_{j=1}^M z_j$  and  $R = (r_1, \dots, r_M)$ . The Schrödinger operator of such a molecule is (12) with  $V(x)$  given by

$$V(x) = - \sum_{j=1}^M \frac{z_j}{|x - r_j|}, \quad (14)$$

the potential of the interaction of an electron with the nuclei. Denote by  $E(Z, R)$  the ground state energy of such a molecule. One of the most elementary questions here is to understand asymptotic behaviour of  $E(Z, R)$  as  $\sum z_j$  tends to  $\infty$ . To make this problem physically and mathematically interesting we have also to scale  $r_j$ . Knowing such asymptotic behaviour is one of the key inputs in the proof of the previous theorem.

**TF Gas.** L. H. Thomas and E. Fermi have suggested in 1927 that a large Coulomb system (atom or molecule) in the ground state (= eigenfunction corresponding to

the lowest eigenvalue) looks like a classical gas but with Pauli principle, namely, there could be at most 2 electrons per volume  $(2\pi)^3$  in the phase space. Such an object is called now the *Thomas-Fermi gas*. Its states are described by the *electron density*  $\varrho \geq 0$  on  $\mathbb{R}^3$  normalized as

$$\int \varrho = N = \# \text{ of electrons}.$$

The energy of the Thomas-Fermi gas is given by a simple non-linear functional

$$\mathcal{E}^{\text{TF}}(\varrho) = \gamma \int \varrho^{5/3} + \int V\varrho + \frac{1}{2} \int \varrho(|x|^{-1} * \varrho), \quad (15)$$

where  $\gamma = (3/5)(3\pi^2)^{2/3}$  and  $V(x)$  is given by (14). The ground state energy,  $E^{\text{TF}}(Z, R)$ , is the infimum (in fact, minimum if  $N \leq \sum z_j$ ) of this functional. It has the following scaling property

$$E^{\text{TF}}(Z, R) = a^{7/3} E^{\text{TF}}(a^{-1}Z, a^{1/3}R).$$

Hence

$$E^{\text{TF}}(Z, R) = O(|Z|^{7/3}).$$

**Asymptotics.** The next theorem shows that the Thomas-Fermi theory is asymptotically correct for large  $Z$  systems but only to the leading order.

**Theorem.** *Let  $Z \rightarrow \infty$  along a given direction and let the mutual distances between  $r_j$  be bounded from below by  $\text{const}|Z|^{-2/3+\epsilon}$ . Then*

$$E(Z, R) = E^{\text{TF}}(Z, R) + \frac{1}{4} \sum z_j^2 + o(|Z|^2). \quad (16)$$

It will follow from the analysis below that the leading term on the r.h.s. represents the quasiclassical energy of the bulk of electrons and the second term, the quantum spectrum of Coulomb singularities.

The leading term in (16) was obtained in [LiebSim 1977]. The second term of asymptotics was conjectured by J. M. C. Scott in 1952 as a contribution of those electrons which move very close to the nuclei (see [LiebSim 1977 and Lieb 1981] for a discussion). For atoms the Scott conjecture was proven in [Hughes 1990, SiedWeik 1987, 1989] and for molecules, in [IvSig 1990]. A proof of the next,  $z^{5/3}$  term for atoms is announced in [FeffSeco 1990b]. The approach in [Hughes 1990, SiedWiek 1987, 1989, FeffSeco 1990b] is based on an expansion in angular momentum channels. This is possible since the electron interaction with the nucleus  $V(x)$ , is spherically symmetric in atoms. The problem is then reduced to a one-dimensional one which is treated by the standard WKB method. The proof in [IvSig 1990] is rather general. I will make a few comments about it.

**Mean Field Theory.** The main utility of the Thomas-Fermi theory for us is to provide a sufficiently simple Schrödinger operator approximating the rather complex  $H$ . Let  $\psi$  be the ground state of  $H$ , i.e. the eigenfunction corresponding to the smallest eigenvalue. Consider the random variable of electron density  $\sum_{i=1}^N \delta(x - x_i)$  with the probability distribution  $|\psi(x_1, \dots, x_N)|^2 dx_1, \dots, dx_N$ . The

ideas of Thomas and Fermi suggest that for large  $N$  this random variable is close to some mean electron density, namely the one,  $\varrho^{\text{TF}}(x)$ , which minimizes  $\mathcal{E}^{\text{TF}}(\varrho)$ . Hence the potential experienced by any one electron is approximately

$$\phi(x) = V(x) + |x|^{-1} * \varrho^{\text{TF}}(x),$$

i.e. the one produced by the nuclei screened by this mean electron density  $\varrho^{\text{TF}}$ . Thus we introduce the mean-field (or quasiparticle) Schrödinger operator

$$H^{\text{ind}} = \sum_{i=1}^N (-\Delta_i + \phi(x_i)) - D$$

acting on  $\bigwedge_{i=1}^N L^2(\mathbb{R}^3 \times \mathbb{Z}_2)$ , where  $D$  is a number compensating for overcounting the electron-electron interaction in  $\phi$ ,

$$D = \frac{1}{2} \iint \frac{\varrho^{\text{TF}}(x)\varrho^{\text{TF}}(y)}{|x-y|} dx dy.$$

This operator describes independent (quasi-) electrons moving in an external potential  $\phi(x)$ . It was realized for some time that in many calculations  $H$  can be replaced by  $H^{\text{ind}}$ .

Some potential theory, Lieb-Thirring inequality, which combines uncertainty and Pauli principles, and a simple version of quasiclassical estimates discussed below yield the following result (cf. Lieb 1981, Hughes 1990, SiedWeik 1987, 1990, FeffSeco 1990b).

**Theorem.** *Let  $E^{\text{ind}}(Z, R)$  be the ground state energy of  $H^{\text{ind}}$ . Then*

$$E(Z, R) = E^{\text{ind}}(Z, R) + O(|Z|^{5/3}).$$

**The One-Body Problem.** The separation of variables on  $\bigwedge_{i=1}^N L^2(\mathbb{R}^3 \times \mathbb{Z}_2)$  gives

$$E^{\text{ind}}(Z, R) = \sum_{i=1}^N E_i - D_{\text{TF}},$$

where  $E_1, E_2, \dots$  are the eigenvalues of the one-particle operator  $P = -\Delta + \phi(x)$  acting on  $L^2(\mathbb{R}^3 \times \mathbb{Z}_2)$ , labeled in order of their increase and counting their multiplicities. This is a much used relation in Quantum Physics and is a consequence of the Pauli principle: at most two electrons (the double degeneracy corresponding to  $\mathbb{Z}_2$ ) per a quantum state. We show

**Theorem.** *Let  $Z \rightarrow \infty$  along a given direction and let the mutual distances between the  $r_j$ 's  $\geq \text{const}|Z|^{-2/3+\epsilon}$ . Then*

$$\sum_{i=1}^N E_i = \text{Weyl} + \text{Scott} + o(|Z|^2), \quad (17)$$

where, with  $p(x, \xi) = |\xi|^2 + \phi(x)$ , the symbol of  $P$ ,

$$\text{Weyl} = \iint_{p \leq 0} p dx d\xi$$

and

$$\text{Scott} = \frac{1}{8} \Sigma z_j^2.$$

The last two theorems yield the main result.

**Quasiclassical Asymptotics.** Now we discuss the proof of (17). Let  $\lambda$  be a fixed vector in  $\mathbb{R}^M$  with non-negative components and let  $Z = \beta^{-3}\lambda$ . The scaling  $x \rightarrow \beta x$  and properties of  $\phi(x)$  show that the large  $Z$  problem is, in fact, the quasiclassical problem with  $\beta = O(|Z|^{-1/3})$  playing the role of the Planck constant. Thus we have to find the quasiclassical asymptotic (as  $\beta \rightarrow 0$ ) for the sum of the first  $N = |Z|$  eigenvalues of the *one-particle* Schrödinger operator  $P$ . The problem is that  $\phi(x)$  is singular. Thus standard quasiclassical methods based on pseudodifferential calculus do not work here. Our method originates in ideas of [Ivr 1986].

There are two ingredients in our proof. First of all we estimate global quantities through local ones. For instance, we study

$$\text{tr}(\psi(x)g(P)), \quad (18)$$

where  $g(\lambda) = \lambda$  for  $\lambda \leq 0$  and  $= 0$  for  $\lambda \geq 0$ . If  $\psi \equiv 1$ , then the trace above is just the sum of negative eigenvalues of  $P$ . We take for  $\psi$  smooth functions localized outside of the singularities of the potentials. Then it is not difficult to obtain asymptotic expansion in the quasiclassical parameter  $\beta$  of the trace (18). Pseudodifferential calculus provides convenient tools for such a purpose. Adapting a standard technique, one represents  $g(P)$  as

$$g(P) = \int \hat{g}(t)e^{-iPt}dt,$$

where  $\hat{g}(t)$  is the Fourier transform of  $g$ . The evolution operator  $e^{-iPt}$  is then approximated on a small time interval to any power in  $\beta$  by Fourier integral operators in the spirit of the geometrical optics. Such an approximation is possible because of finite speed of propagation of singularities for the Schrödinger equation: for sufficiently small times and for bounded energies the singularities of  $\phi(x)$  do not reach  $\text{supp } \psi$ . The appropriate Fourier integral operators are then expanded by the method of stationary phase. The information about (18) is recovered using the Tauberian technique. However, the remainder estimates here depend on  $\text{supp } \psi$  and on estimates of  $\partial^\nu \phi$  on  $\text{supp } \psi$ .

The second ingredient is a multiscale analysis. There are three scales in the problem: momentum scale determined by the quasiclassical parameter  $\beta = O(|Z|^{-1/3})$ , space scale,  $l(x)$ , determined by how the potential differentiates and the energy scale,  $f(x)$ , determined by the size of the potential. The first scale is constant while the other two depend on  $x$ . In our problem

$$l(x) = \text{dist of } x \text{ to the singularities}$$

and  $f(x) = l(x)^{-1}$ . At each point outside of the singularities we rescale the problem using the scales at this point in such a way that the problem is mapped

into a model one, i.e. the one with a potential  $U$  obeying  $|\partial^\alpha U(x)| \leq C_\alpha$  on a unit ball, with the effective quasiclassical parameter

$$\alpha_{\text{eff}}(x) = \frac{\beta}{l(x)f(x)^{1/2}},$$

which depends on all the scales. The new problem admits a quasiclassical expansion discussed above with a remainder bound independent of the singular structure of  $\phi(x)$ . This implies an expansion for the original problem outside of small balls around the singularities of  $\phi(x)$ . Inside each of those balls we analyze the problem differently. Namely, we replace the potential by its leading term near the singularity and solve the quasiclassical problem for this truncated operator more accurately. Leading terms near and outside the singularities combine into a single Weyl term over  $\mathbb{R}^3$  which gives the Thomas-Fermi energy. Precise quasiclassical expansion of low-lying eigenvalues of the truncated problem near the singularities yields the Scott correction.

## 4. Conclusion

The rigorous Quantum Mechanics of many-particle systems is a fast developing branch of Mathematical Physics. This paper is not a comprehensive review of the subject. A rather versatile account of the Schrödinger operators as well as many references can be found in [CFKS]. What we attempted here is a brief glimpse into some of the recent developments and trends, just a small sample of many fresh and exciting problems with which this field abounds.

*Acknowledgements.* This paper was written while the author was visiting the Institut für Theoretische Physik, ETH-Zürich; the author is grateful to J. Fröhlich and W. Hunziker for the hospitality. It is a pleasure to thank V. Bach, J. Fröhlich, W. Hunziker, B. Simon and J.-P. Solovej for reading the manuscript and making many useful remarks. This paper sums up results of the joyful collaboration with A. Soffer, V. Ivrii, J.-P. Solovej and L. Seco.

## References

- Bach, V. (1991): Ionization energies of bosonic Coulomb systems. *Lett. Math. Phys.* **21**, 139–149
- Baumgartner, B. (1984): On the Thomas-Fermi-von-Weizsäcker and Hartree energies as functions of the degree of ionization. *J. Phys.* **A17**, 1593–1602
- Benguria, R., Lieb, E.H. (1983): Proof of the stability of highly negative ions in the absence of the Pauli principle. *Phys. Rev. Lett.* **50**, 1771–1774
- Cycon, H.L., Froese, R.G., Kirsch, W., Simon, B. (1987): Schrödinger operators. Springer, Berlin Heidelberg New York
- Enss, V. (1985): Quantum scattering theory for two- and three-body systems with potentials of short- and long-range. In: Schrödinger Operators (S. Graffi, ed.). (Lecture Notes in Mathematics, vol. 1159.) Springer, Berlin Heidelberg New York
- Fefferman, C., Seco, L. (1990a): Asymptotic neutrality of large ions. *Comm. Math. Phys.* **128**, 109–130
- Fefferman, C., Seco, L. (1990b): On the energy of a large atom. *Bulletin AMS* **23**, 525–530

- Graf, G.-M. (1990): Asymptotic completeness for  $N$ -body short-range quantum systems: a new proof. *Comm. Math. Phys.* **132**, 73–101
- Hughes, W. (1990): An atomic energy lower bound that gives Scott's correction. *Adv. Math.* **79**, 213
- Ivrii, V. (1986): Weyl's asymptotic for the Laplace-Beltrami operator in Riemann polyhedra. *Dokl. Akad. Nauk SSSR* **38**, 35–38
- Ivrii, V., Sigal, I.M. (1990): Asymptotics on the ground state energies of large Coulomb systems. *Ann. Math.* (submitted)
- Kitada, H. (1990): On the completeness of  $N$ -body wave operators. Preprint, Tokyo
- Lieb, E. (1981): Thomas-Fermi and related theories of atoms and molecules. *Rev. Modern Phys.* **53**
- Lieb, E., Sigal, I.M., Simon, B., Thirring, W. (1988): Approximate neutrality of large  $Z$ -ions. *Comm. Math. Phys.* **116**, 635–644
- Lieb, E., Simon, B. (1977): Thomas-Fermi theory of atoms, molecules and solids. *Adv. Math.* **23**, 22–116
- Mourre, E. (1981): Absence of singular continuous spectrum of certain self-adjoint operators. *Comm. Math. Phys.* **78**, 391–408
- Perry, P., Sigal, I.M., Simon, B. (1981): Spectral analysis of  $N$ -body Schrödinger operators. *Ann. Math.* **114**, 519–567
- Ruskai, M.B. (1982): Absence of discrete spectrum of highly negative ions. *Comm. Math. Phys.* **82**, 457–469
- Seco, L., Sigal, I.M., Solovej, J.-P. (1990): Bound on the ionization energy of large atoms. *Comm. Math. Phys.* **131**, 307–315
- Siedentop, H., Weikard, R. (1987): On the leading energy correction for the statistical model of the atom: interacting case. *Comm. Math. Phys.* **112**, 471–490
- Siedentop, H., Weikard, R. (1989): On the leading correction of the Thomas-Fermi model: lower bound. *Invent. math.* **97**, 159–193
- Sigal, I.M. (1982): Geometric methods in the quantum many-body problem. Nonexistence of very negative ions. *Comm. Math. Phys.* **85**, 309–324
- Sigal, I.M. (1989): Geometrical theory of resonances in multi-particle systems. In: *Proceedings IX Congress Math. Phys.* (B. Simon et al., eds.). Adam Hilger
- Sigal, I.M., Soffer, A. (1987): The  $N$ -particle scattering problem: asymptotic completeness for short-range systems. *Ann. Math.* **125**, 35–108
- Sigal, I.M., Soffer, A. (1990a): Long-range many-body scattering. Asymptotic clustering for Coulomb-type potential. *Invent. math.* **99**, 115–143
- Sigal, I.M., Soffer, A. (1990b): Asymptotic completeness for 4 and 5 particle systems with the Coulomb-type potentials. University of Toronto. Preprint
- Solovej, J.-P. (1990): Asymptotics for bosonic atoms. *Lett. Math. Phys.* **20**, 165–172
- Tamura, H. (1990): Asymptotic completeness for  $N$ -body Schrödinger operators ... Ibaraki. Preprint



# Moduli of Stable Curves, Conformal Field Theory and Affine Lie Algebras

Akihiro Tsuchiya

Department of Mathematics, School of Science, Nagoya University, Nagoya 464-01, Japan

## §1. Introduction

Conformal field theory was initiated by Belavin, Polyakov and Zamolodchikov [1] as 2-dimensional quantum field theory describing the 2-dimensional critical phenomena. Conformal field theory is characterized by infinite-dimensional symmetries such as Virasoro algebra, and its correlation functions are characterized by differential equations arising from representations of infinite-dimensional Lie algebras.

In this article, we report our works with Y. Kanie [12], K. Ueno and Y. Yamada [13] concerning the construction of conformal field theory on the universal family of stable curves under the gauge symmetries associated with integrable representations of Lie algebras.

We realize it by constructing a coherent  $\mathcal{O}_{M_{g,N}^{(1)}}$ -module  $\mathcal{V}_A$  over the modular stack  $M_{g,N}^{(1)}$  of  $N$ -pointed stable curves of genus  $g$  with first order infinitesimal structure, the sheaf of twisted first order differential operators  $\mathcal{D}_{M_{g,N}^{(1)}}^1(-\log D_{g,N}^{(1)} : c_\nu)$ , which is the geometric counterpart of the Virasoro algebra, and the action of  $\mathcal{D}_{M_{g,N}^{(1)}}^1(-\log D_{g,N}^{(1)} : c_\nu)$  on  $\mathcal{V}_A$ . The solution sheaf of  $\mathcal{V}_A$  gives what physicists call the current-conformal block which is the most fundamental object in conformal field theory. Moreover it turns out to be locally free and the factorization property at the normal crossing divisors  $D_{g,N}^{(1)}$  holds.

The monodromy of these solution sheaves give the representations of the central extension  $\widehat{\Gamma}_{g,N}$  of the mapping class group  $\Gamma_{g,N}$ , by the multiplicative subgroup  $K$  of  $\mathbf{C}^*$  generated by  $e^{2\pi\sqrt{-1}c_\nu/24}$ .

The construction of conformal blocks was also done using a quite different method (topological point of view) by G. Moore and N. Seiberg [9].

## §2. Integrable Modules of Affine Lie Algebras and the Sugawara Form

Let  $\mathfrak{g}$  be a simple Lie algebra over  $\mathbf{C}$ . We fix a Cartan subalgebra  $\mathfrak{h}$ , a simple root system  $\Sigma$ , and the invariant bilinear form  $( , )$  on  $\mathfrak{g}$  normalized by  $(\theta, \theta) = 2$

for the highest root  $\theta$ . Let  $P_+$  denote the set of dominant integral weights of  $\mathfrak{g}$  for  $\lambda \in P_\ell$ . We denote by  $V_\lambda$  the irreducible  $\mathfrak{g}$ -module with highest weight  $\lambda$ . And  $V_0$  is the one dimensional trivial  $\mathfrak{g}$ -module,  $V_{\lambda^*}$  denotes the contragradient  $\mathfrak{g}$ -module. For each positive integer  $\ell$ , we define the finite subset  $P_\ell$  of  $P_+$  by  $\{\lambda \in P_\ell; 0 \leq (\theta, \lambda) \leq \ell\}$ . We fix an orthonormal basis  $\{J_a\}$  of  $\mathfrak{g}$ .

Let  $\widehat{\mathfrak{g}} = \mathfrak{g} \otimes \mathbf{C}((\xi)) \oplus \mathbf{C}\epsilon$  denote the associated affine Lie algebra of  $\mathfrak{g}$ . For  $X \in \mathfrak{g}$ ,  $f \in \mathbf{C}((\xi))$  and  $n \in \mathbb{Z}$ , we set  $X[f] = X \otimes f$  and  $X(n) = X \otimes \xi^n$ .

In the sequel we fix a positive integer  $\ell$ . The set  $P_\ell$  parametrizes the integrable highest weight  $\widehat{\mathfrak{g}}$ -modules of level  $\ell$ , and for each  $\lambda \in P_\ell$  we denote the associated integrable  $\widehat{\mathfrak{g}}$ -module by  $\mathcal{H}_\lambda$ . For each  $\lambda \in P_\ell$ , put  $\Delta_\lambda = (\lambda, \lambda + 2\varrho)/2(\ell + g^*)$ , where  $\varrho$  is half the sum of the positive roots of  $\mathfrak{g}$ .

For each  $n \in \mathbb{Z}$ , we define the Sugawara operator  $T(n)$  acting on  $\mathcal{H}_\lambda$  by

$$T(n) = \frac{1}{2(\ell + g^*)} \sum_a \sum_{k \in \mathbb{Z}} \circ J_a(k) J_a(n-k) \circ \quad (2.1)$$

where the symbol  $\circ \quad \circ$  denotes the so-called normal ordering and  $g^*$  denotes the dual Coxeter number of  $\mathfrak{g}$ . For an element  $\ell = \sum_n b_n \xi^{n+1} d/d\xi \in \mathbf{C}((\xi))d/d\xi$ , define an operator  $T[\ell]$  on  $\mathcal{H}_\lambda$  by  $T[\ell] = -\sum_n b_n T(n)$ . Then we have the following fundamental relations of the operators on  $\mathcal{H}_\lambda$ .

$$\begin{aligned} [X[f], Y[g]] &= [X, Y][fg] + \ell(X, Y) \operatorname{Res}(df \cdot g) \operatorname{id} \\ [T[\ell], X[f]] &= X[\ell(f)] \\ [T[\ell_1], T[\ell_2]] &= T[[\ell_1, \ell_2]] + \frac{c_v}{12} \operatorname{Res} \frac{d^3 \ell_1}{d\xi^3} \ell_2 d\xi \operatorname{id} \end{aligned} \quad (2.2)$$

where  $c_v = \ell \dim \mathfrak{g}/(\ell + g^*)$ .

Consider the automorphism group  $\mathcal{D}$  of the  $\mathbf{C}$ -algebra  $\mathbf{C}[[\xi]]$ . Then an element  $h$  of  $\mathcal{D}$  is represented by a formal power series  $h(\xi) = a_0 \xi + a_1 \xi^2 + \dots$ ,  $a_0 \neq 0$ . For each  $n \geq 0$ , define the normal subgroup  $\mathcal{D}_n$  of  $\mathcal{D}$  by  $\mathcal{D}_n = \{h(\xi) = \xi + a_n \xi^{n+1} + \dots\}$ . Since each element  $h$  of  $\mathcal{D}_1$  is uniquely written as  $h = \exp(\ell)$ ,  $\ell \in \mathbf{C}[[\xi]]\xi^2 d/d\xi$ , we define the operator  $G[h]$  of  $\mathcal{H}_\lambda$  by  $G[h] = \exp(T[\ell])$ . These operators  $G[h]$  define a representation of  $\mathcal{D}_1$  on  $\mathcal{H}_\lambda$ .

### §3. Local Universal Family of $N$ -Pointed Stable Curves of Genus $g$

Let  $(g, N)$  be a pair of non-negative integers with  $2g - 2 + N \geq 1$ . Consider a local universal family of  $N$ -pointed connected stable curves of genus  $g$ .

$$\mathcal{F} = (\pi : C \rightarrow B : s_1, \dots, s_N) \quad (3.1)$$

where  $s_j : B \rightarrow C$  are cross-sections of  $\pi$ , and set  $S_j = s_j(B)$  and  $S = \bigcup_j S_j$ .

Let  $D$  be the discriminant locus of  $\pi : C \rightarrow B$ . Then  $D$  is a normal crossing divisor of  $B$ , and  $B$  is a  $3g - 3 + N$  dimensional complex manifold.

For each non-negative integer  $n$  or  $n = \infty$  we consider the associated local universal family of  $N$ -pointed stable curves of genus  $g$  with  $n$ -th infinitesimal neighborhoods.

$$\mathcal{F}^{(n)} = (\pi^{(n)} : C^{(n)} \longrightarrow B^{(n)} : s_1^{(n)}, \dots, s_N^{(n)}; t_1^{(n)}, \dots, t_N^{(n)}) \quad (3.2)$$

where  $t_j^{(n)} : \mathcal{O}_{B^{(n)}}[[\zeta]]/(\zeta^{n+1}) \rightarrow \mathcal{O}_{C^{(n)}}/I_{S_j^{(n)}}^{n+1}$  are  $\mathcal{O}_{B^{(n)}}$ -algebra isomorphisms.

We get the sequence of fibering

$$B = B^{(0)} \leftarrow B^{(1)} \leftarrow \cdots \leftarrow B^{(\infty)} \quad (3.3)$$

where  $p_m^n = p : B^{(n)} \rightarrow B^{(m)}$  is a principal fibering with structure group  $\mathcal{D}_m^{n \oplus N}$ , where  $\mathcal{D}_m^n = \mathcal{D}_m / \mathcal{D}_n$ .

By the local universality of the system  $\mathcal{F}^{(n)}$ , we have the following isomorphism of  $\mathcal{O}_{B^{(n)}}$ -modules.

$$\varrho_n : \mathcal{O}_{B^{(n)}}(-\log D^{(n)}) \xrightarrow{\sim} R^1\pi_*^{(n)}(\mathcal{O}_{C^{(n)}/B^{(n)}}(-(n+1)S^{(n)})) \quad (3.4)$$

where  $D^{(n)} = p^{-1}(D)$ ,  $S_j^{(n)} = s_j^{(n)}(B^{(n)})$ ,  $S^{(n)} = \cup S_j^{(n)}$ .

With each local universal family  $\mathcal{F} = (\pi : C \rightarrow B : s_1, \dots, s_N)$ , we associate the following  $\mathcal{O}_{B^{(n)}}$ -modules

$$\begin{aligned} K_{\widehat{S}_j^{(n)}} &= \varprojlim_m \mathcal{O}_{C^{(n)}}(*S_j^{(n)})/I_{S_j^{(n)}}^{m+1} \\ \mathcal{O}_{\widehat{S}_j^{(n)}/B^{(n)}}(*) &= \mathcal{D}\text{er}_{\mathcal{O}_{B^{(n)}}}(K_{\widehat{S}_j^{(n)}}, K_{\widehat{S}_j^{(n)}}). \end{aligned} \quad (3.5)$$

The  $\mathcal{O}_{B^{(n)}}$ -module  $\mathcal{O}_{\widehat{S}_j^{(n)}/B^{(n)}}(*)$  has the canonical Lie algebra structure over  $\mathcal{O}_{B^{(n)}}$ . We denote its Lie bracket by  $[ , ]_0$ .

For  $n = \infty$  we have the following canonical  $\mathcal{O}_{B^{(\infty)}}$ -module isomorphisms.

$$\begin{aligned} \tilde{t}_j^{(\infty)} : \mathcal{O}_{B^{(\infty)}}((\xi_j)) &\xrightarrow{\sim} K_{\widehat{S}_j^{(\infty)}} \\ \tilde{t}_j^{(\infty)} : \mathcal{O}_{B^{(\infty)}}((\xi_j)) \frac{d}{d\xi_j} &\xrightarrow{\sim} \mathcal{O}_{\widehat{S}_j^{(\infty)}/B^{(\infty)}}(*). \end{aligned} \quad (3.6)$$

We have the following canonical inclusion mappings.

$$\begin{aligned} \pi_*^{(n)} \mathcal{O}_{C^{(n)}}(*S^{(n)}) &\hookrightarrow \sum_{j=1}^N K_{\widehat{S}_j^{(n)}} \\ \pi_*^{(n)} \mathcal{O}_{C^{(n)}/B^{(n)}}(*S^{(n)}) &\hookrightarrow \sum_{j=1}^N \mathcal{O}_{\widehat{S}_j^{(n)}/B^{(n)}}(*). \end{aligned} \quad (3.7)$$

Now consider the following condition (Q) for the system  $\mathcal{F} = (\pi : C \rightarrow B : s_1, \dots, s_N)$ .

(Q) For each point  $b \in B$ , and for each irreducible component  $C$  of  $C_b = \pi^{-1}(b)$ , there exists  $j$  with  $s_j(b) \in C$ .

**Proposition 3.1.** *Under the condition (Q) for  $\mathcal{F}$ , we have the following surjective  $\mathcal{O}_{B^{(n)}}$ -module homomorphism*

$$\theta : \sum_{j=1}^N \Theta_{\widehat{S}_j^{(n)}/B^{(n)}}(*) \longrightarrow \Theta_{B^{(n)}}(-\log D^{(n)}). \quad (3.8)$$

Furthermore, we introduce the Lie bracket  $[ , ]$  on  $\sum_{j=1}^N \Theta_{\widehat{S}_j^{(n)}/B^{(n)}}(*)$  by

$$[v_1, v_2] = [v_1, v_2]_0 + \theta(v_1)(v_2) - \theta(v_2)(v_1). \quad (3.9)$$

Then the map  $\theta$  is a homomorphism of Lie algebras.

As the group  $\mathcal{D}^{\oplus N}$  acts on  $B^{(\infty)}$ , the induced action of  $h = (h^1, \dots, h^N) \in \mathcal{D}^{\oplus N}$  on  $f = \sum_{j=1}^N \sum_n a_n^j \xi_j^n \in \sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j))$  and  $\ell = \sum_{j=1}^N \sum_n b_n^j \xi_j^{n+1} d/d\xi_j \in \sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j)) d/d\xi_j$  are defined by  $\pi(h)(f) = \sum_j \sum_n \varrho(h)(a_n^j) h^j(\xi_j^{n+1})$  and  $\pi(h)(\ell) = \sum_j \sum_n \varrho(h)(b_n^j) \text{Ad}(h^j)(\xi_j^{n+1} d/d\xi_j)$  respectively.

**Proposition 3.4.** *The action of  $\mathcal{D}^{\oplus N}$  preserves the subspace  $\pi_*^{(\infty)} \mathcal{O}_{B^{(\infty)}}(*S^{(\infty)})$  of  $\sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j))$ .*

We remark that the invariant part of the action of  $\mathcal{D}_n^{\oplus N}$  on  $\sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j))$  and  $\sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j)) \frac{d}{d\xi_j}$  are  $\sum_{j=1}^N K_{\widehat{S}_j^{(n)}}$  and  $\sum_{j=1}^N \Theta_{\widehat{S}_j^{(n)}/B^{(n)}}(*)$  respectively.

## §4. Sheafification and the Module $\mathcal{V}_{\lambda}$

In this section, we assume the condition (Q) for any local universal family  $\mathcal{F}$ .

With each system  $\mathcal{F} = (\pi : C \rightarrow B : s_1, \dots, s_N)$ , we associate the  $\mathcal{O}_{B^{(\infty)}}$ -Lie algebra sheaves

$$\begin{aligned} \widehat{\mathfrak{g}}(\mathcal{F}^{(n)}) &= \sum_{j=1}^n \mathfrak{g} \otimes K_{\widehat{S}_j^{(n)}} \oplus \mathcal{O}_{B^{(n)}} c \\ \widehat{\mathfrak{g}}_-(\mathcal{F}^{(n)}) &= \mathfrak{g} \otimes \pi_*^{(n)} \mathcal{O}_{C^{(n)}}(*S^{(n)}) \subseteq \widehat{\mathfrak{g}}(\mathcal{F}^{(n)}). \end{aligned} \quad (4.1)$$

The Lie algebra structure is given by

$$\left[ \sum_{j=1}^n X_j \otimes f_j, \sum_{j=1}^n Y_j \otimes g_j \right] = \sum_{j=1}^n [X_j, Y_j] \otimes f_j g_j + \sum_{j=1}^n \text{Res}(g_j d f_j)(X_j, Y_j) c. \quad (4.2)$$

Then  $\widehat{\mathfrak{g}}_-(\mathcal{F}^{(n)})$  is a Lie subalgebra of  $\widehat{\mathfrak{g}}(\mathcal{F}^{(n)})$ . In the case of  $n = \infty$ , we have

$$\widehat{\mathfrak{g}}(\mathcal{F}^{(\infty)}) = \sum_{j=1}^N \mathfrak{g} \otimes \mathcal{O}_{B^{(\infty)}}((\xi_j)) \oplus \mathcal{O}_{B^{(\infty)}} c. \quad (4.3)$$

The Lie group  $\mathcal{D}^{\oplus N}$  acts on  $\widehat{\mathfrak{g}}(\mathcal{F}^{(\infty)})$  as Lie algebra automorphisms, and the sub-algebra  $\widehat{\mathfrak{g}}_-(\mathcal{F}^{(\infty)})$  is preserved by the action of  $\mathcal{D}^{\oplus N}$ .

For each  $\lambda = (\lambda_1, \dots, \lambda_N)$ , a quasi-coherent  $\mathcal{O}_{B^{(\infty)}}$ -module  $\mathcal{H}_\lambda(\mathcal{F}^{(\infty)})$  is defined by

$$\mathcal{H}_\lambda(\mathcal{F}^{(\infty)}) = \mathcal{O}_{B^{(\infty)}} \otimes \mathcal{H}_\lambda, \quad \mathcal{H}_\lambda = \mathcal{H}_{\lambda_1} \otimes \cdots \otimes \mathcal{H}_{\lambda_N}. \quad (4.4)$$

Then the Lie algebra sheaf  $\widehat{\mathfrak{g}}(\mathcal{F}^{(\infty)})$  acts on  $\mathcal{H}_\lambda(\mathcal{F}^{(\infty)})$   $\mathcal{O}_{B^{(\infty)}}$ -linearly. Define the  $\mathcal{O}_{B^{(\infty)}}$ -modules  $\mathcal{H}'_\lambda(\mathcal{F}^{(\infty)})$  and  $\mathcal{V}'_\lambda(\mathcal{F}^{(\infty)})$  by

$$\begin{aligned} \mathcal{H}'_\lambda(\mathcal{F}^{(\infty)}) &= \widehat{\mathfrak{g}}_-(\mathcal{F}^{(\infty)}) \cdot \mathcal{H}_\lambda(\mathcal{F}^{(\infty)}) \\ \mathcal{V}'_\lambda(\mathcal{F}^{(\infty)}) &= \mathcal{H}_\lambda(\mathcal{F}^{(\infty)}) / \mathcal{H}'_\lambda(\mathcal{F}^{(\infty)}). \end{aligned} \quad (4.5)$$

Since the Lie group  $\mathcal{D}_1^{\oplus N}$  acts on  $\mathcal{H}_\lambda(\mathcal{F}^{(\infty)})$ , preserving  $\mathcal{H}'_\lambda(\mathcal{F}^{(\infty)})$ , the  $\mathcal{O}_{B^{(\infty)}}$ -module  $\mathcal{V}'_\lambda(\mathcal{F}^{(\infty)})$  has the structure of a  $\mathcal{D}_1^{\oplus N}$ -module.

Define the  $\mathcal{O}_{B^{(1)}}$ -modules  $\mathcal{H}_\lambda(\mathcal{F}^{(1)})$ ,  $\mathcal{H}'_\lambda(\mathcal{F}^{(1)})$  and  $\mathcal{V}'_\lambda(\mathcal{F}^{(1)})$  as the invariant part of the  $\mathcal{D}_1^{\oplus N}$ -actions on  $\mathcal{H}_\lambda(\mathcal{F}^{(\infty)})$ ,  $\mathcal{H}'_\lambda(\mathcal{F}^{(\infty)})$  and  $\mathcal{V}'_\lambda(\mathcal{F}^{(\infty)})$  respectively. Then the  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{H}_\lambda(\mathcal{F}^{(1)})$  has a canonical structure of  $\widehat{\mathfrak{g}}(\mathcal{F}^{(1)})$ -module and we have the following canonical  $\mathcal{O}_{B^{(1)}}$ -module isomorphism.

$$\begin{aligned} \mathcal{H}'_\lambda(\mathcal{F}^{(1)}) &\simeq \widehat{\mathfrak{g}}_-(\mathcal{F}^{(1)}) \mathcal{H}_\lambda(\mathcal{F}^{(1)}) \\ \mathcal{V}'_\lambda(\mathcal{F}^{(1)}) &\simeq \mathcal{H}_\lambda(\mathcal{F}^{(1)}) / \mathcal{H}'_\lambda(\mathcal{F}^{(1)}). \end{aligned} \quad (4.6)$$

The  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{V}'_\lambda(\mathcal{F}^{(1)})$  is the most fundamental object in conformal field theory. Our first main theorem is the following one.

**Theorem I.** *The module  $\mathcal{V}'_\lambda(\mathcal{F}^{(1)})$  is a coherent  $\mathcal{O}_{B^{(1)}}$ -module.*

For the proof of this theorem, we use Gabber's theorem on the involutiveness of the characteristic variety [Gabber 5].

## §5. Differential Equations Defining Conformal Blocks

In this section we also assume the condition (Q), unless otherwise stated. Now to each  $\mathcal{F}$ , we associate the following sheaf version of the Virasoro algebra.

$$\text{Vir}(\mathcal{F}^{(\infty)} : c_v) = \sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j)) \frac{d}{d\xi_j} \oplus \mathcal{O}_{B^{(\infty)}} \quad (5.1)$$

with the Lie bracket

$$[(\ell, r), (m, s)] = \left( [\ell, m], \frac{c_v}{12} \sum_{j=1}^N \text{Res}_{\xi_j=0} \left( \frac{d^3 \ell_j}{d \xi_j^3} m_j d \xi_j \right) + \theta(\ell)(s) - \theta(m)(r) \right) \quad (5.2)$$

where

$$\begin{aligned}\ell &= \left( \ell_1 \frac{d}{d\xi_1}, \dots, \ell_N \frac{d}{d\xi_N} \right), \quad m = \left( m_1 \frac{d}{d\xi_1}, \dots, m_N \frac{d}{d\xi_N} \right) \\ &\in \sum_{j=1}^N \mathcal{O}_{B^{(\infty)}}((\xi_j)) \frac{d}{d\xi_j} \text{ and } r, s \in \mathcal{O}_{B^{(\infty)}}.\end{aligned}$$

The group  $\mathcal{D}^{\oplus N}$  acts on  $\text{Vir}(\mathcal{F}^{(\infty)} : c_v)$  as Lie algebra automorphisms as follows. For  $h = (h_1, \dots, h_N) \in \mathcal{D}^{\oplus N}$ ,  $v = (\ell, r) \in \text{Vir}(\mathcal{F}^{(\infty)}, c_v)$ ,  $\ell = (\ell_1 d/d\xi_1, \dots, \ell_N d/d\xi_N)$ ,

$$\begin{aligned}\pi(h)(\ell, r) &= (\pi(h)(l), \varrho(h)(r)) \\ &\quad + \frac{c_v}{12} \sum_{j=1}^N \text{Res}_{\xi_j=0}(\pi(h)(\ell_j)\{h_j(\xi_j), \xi_j\} \left( \frac{dh_j(\xi_j)}{d\xi_j} \right)^{-1} d\xi_j)),\end{aligned}$$

where  $\{h_j(\xi_j), \xi_j\}$  denotes the Schwarzian derivative. Taking the  $\mathcal{D}_1^{\oplus N}$  invariant part, we get the following exact sequence of Lie algebras.

$$0 \longrightarrow \mathcal{O}_{B^{(1)}} \longrightarrow \text{Vir}(\mathcal{F}^{(1)} : c_v) \longrightarrow \sum_{j=1}^N \Theta_{\widehat{S}_j^{(1)}/B^{(1)}}(*) \longrightarrow 0. \quad (5.3)$$

The Lie algebra sheaf  $\text{Vir}(\mathcal{F}^{(\infty)} : c_v)$  acts on  $\mathcal{H}_{\lambda}(\mathcal{F}^{(\infty)})$  by the following formula

$$D(v)(F \otimes |\Phi\rangle) = \theta(\ell)(F) \otimes \Phi + F \otimes \sum_{j=1}^N \varrho_j(T[\ell_j])|\Phi\rangle + rF \otimes |\Phi\rangle \quad (5.4)$$

where  $F \in \mathcal{O}_{B^{(\infty)}}$ ,  $|\Phi\rangle \in \mathcal{H}_{\lambda}$ ,  $v = (\ell, r) \in \text{Vir}(\mathcal{F}^{(\infty)}, c_v)$ .

**Proposition 5.1.** *The action of  $\text{Vir}(\mathcal{F}^{(\infty)} : c_v)$  on  $\mathcal{H}_{\lambda}(\mathcal{F}^{(\infty)})$  preserves  $\mathcal{H}'_{\lambda}(\mathcal{F}^{(\infty)})$ , so  $\text{Vir}(\mathcal{F}^{(\infty)} : c_v)$  acts on  $\mathcal{V}'_{\lambda}(\mathcal{F}^{(\infty)})$ .*

Finally taking the invariant part of the action  $\mathcal{D}_1^{\oplus N}$ , we get an action of the Lie algebra sheaf  $\text{Vir}(\mathcal{F}^{(1)} : c_v)$  on  $\mathcal{H}_{\lambda}(\mathcal{F}^{(1)})$ ,  $\mathcal{H}'_{\lambda}(\mathcal{F}^{(1)})$  and  $\mathcal{V}'_{\lambda}(\mathcal{F}^{(1)})$ .

**Proposition 5.2.** *For any system  $\mathcal{F}$ , there exists a unique  $\mathcal{O}_{B^{(\infty)}}$ -module homomorphism,*

$$a^{(\infty)} : \pi_*^{(\infty)} \Theta_{C^{(\infty)}/B^{(\infty)}}(*) \oplus \mathcal{O}_{B^{(\infty)}} \longrightarrow \mathcal{O}_{B^{(\infty)}} \quad (5.5)$$

such that for any  $v \in \pi_*^{(\infty)} \Theta_{C^{(\infty)}/B^{(\infty)}}(*) \oplus \mathcal{O}_{B^{(\infty)}}$  and  $|\Phi\rangle \in \mathcal{V}'_{\lambda}(\mathcal{F}^{(\infty)})$  we have,

$$D(v)|\Phi\rangle = a^{(\infty)}(v)|\Phi\rangle. \quad (5.6)$$

This map  $a^{(\infty)}$  is expressed, locally on  $B$ , in the following way. For  $\mathcal{F} = (\pi : C \rightarrow B; s_1, \dots, s_N)$ , taking  $B$  small enough, there exists an element  $\omega \in H^0(C \times_B C, \omega_{C/B}^{\boxtimes 2}(2A))$  such that near the diagonal  $A \subseteq C \times_B C$

$$\omega = \frac{dw dz}{(w-z)^2} + \text{regular at } A, \quad (5.7)$$

where  $w$  and  $z$  represent fiber coordinates of  $\pi : C \rightarrow B$  and  $\omega_{C/B}$  denotes the relative dualizing sheaf of  $\pi$ . Define, then, the associated projective connection by  $S_\omega(z)dz^2 = -6\lim_{w \rightarrow z}\{\omega - dw dz/(w-z)^2\}$ . Then for  $v = ((\ell_1 d/d\xi_1, \dots, \ell_N d/d\xi_N), r)$  we have

$$a^{(\infty)}(v) = \frac{c_v}{12} \sum_{j=1}^N \text{Res}_{\xi_j=0} (\ell_j(\xi_j) S_\omega(\xi_j) d\xi_j) + r. \quad (5.8)$$

Taking the  $\mathcal{D}_1^{\oplus N}$  invariant part we get the  $\mathcal{O}_{B^{(1)}}$ -module homomorphism,

$$a : \pi_*^{(1)} \Theta_{C^{(1)}/B^{(1)}}(*S^{(1)}) \oplus \mathcal{O}_{B^{(1)}} \longrightarrow \mathcal{O}_{B^{(1)}}. \quad (5.9)$$

The sheaf of twisted first-order differential operators on  $B^{(1)}$  is defined by

$$\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v) = \text{Vir}(\mathcal{F}^{(1)} : c_v) / \left( \ker a \oplus \sum_{j=1}^N \Theta_{\widehat{S}_j^{(1)}/B^{(1)}}(-2S_j^{(1)}) \right). \quad (5.10)$$

Then we have the following exact sequence of Lie algebra sheaves on  $B^{(1)}$

$$0 \longrightarrow \mathcal{O}_{B^{(1)}} \longrightarrow \mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v) \longrightarrow \Theta_{B^{(1)}}(-\log D^{(1)}) \longrightarrow 0. \quad (5.11)$$

By Proposition 5.2, we get the following second main theorem of this paper.

**Theorem II.** *On  $\mathcal{V}_\lambda(\mathcal{F}^{(1)})$ , the sheaf of Lie algebras  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  acts as twisted first order differential operators.*

As a corollary of this theorem we have

**Corollary 5.3.** *On  $B^{(1)} - D^{(1)}$ , the coherent  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{V}_\lambda(\mathcal{F}^{(1)})$  is locally free.*

The sheaf of twisted first order differential operators  $\mathcal{D}_{B^{(1)}}^1(-\log D : c_v)$  can be trivialized locally on  $B$  as follows. Take  $\omega \in H^0(C \times_B C, \omega_{C/B}^{\boxtimes 2}(2A))$  satisfying the property (5.7). Then we can associate canonically the  $\mathcal{O}_{B^{(1)}}$ -module isomorphism

$$A_\omega = A_\omega(\mathcal{F}) : \Theta_{B^{(1)}}(-\log D^{(1)} : c_v) \oplus \mathcal{O}_{B^{(1)}} \longrightarrow \mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v) \quad (5.12)$$

which is indeed a Lie algebra isomorphism.

By this isomorphism  $A_\omega(\mathcal{F})$ , the sheaf  $\mathcal{O}_{B^{(1)}}$  has a structure of a  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  module. We denote this  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  module by  $\mathcal{L}_\omega(\mathcal{F}^{(1)})$ . Then we can consider the solution sheaf,

$$\mathcal{H}om_{\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)}(\mathcal{V}_\lambda(\mathcal{F}^{(1)}), \mathcal{L}_\omega(\mathcal{F}^{(1)})). \quad (5.13)$$

A fiber of this sheaf is nothing but what physicists call the space of conformal blocks.

Finally in this section, we remark how to define the sheaves  $\mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)})$  and  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  in the case when the condition (Q) is not satisfied. For any local universal family  $\mathcal{F} = (\pi : C \rightarrow B; s_1, \dots, s_N)$ , not necessarily satisfying the condition (Q), we can take a local universal family of  $(N+k)$  pointed stable curves  $\mathcal{F}' = (\pi' : C' \rightarrow B'; s'_1, \dots, s'_{N+k})$  satisfying the condition (Q) and surjective smooth maps  $F : C' \rightarrow C$  and  $f : B' \rightarrow B$  with following properties:  
1)  $f \circ \pi' = \pi \circ F$   
2)  $F s'_j = s_j f, j = 1, \dots, N$   
3) for any  $b' \in B'$ , put  $b = f(b')$ ,  $C_{b'} = \pi'^{-1}(b')$ ,  $C_b = \pi^{-1}(b)$ . Then, the map  $F_{b'} : (C_{b'} : s'_1(b'), \dots, s'_N(b')) \rightarrow (C_b : s_1(b), \dots, s_N(b))$  is an isomorphism of  $N$ -pointed stable curves. Note that  $D' = f^{-1}(D)$ .

For  $\vec{\lambda} \in P_\ell^N$ , put  $\vec{\lambda} = (\vec{\lambda}, 0, \dots, 0) \in P_\ell^{N+k}$ . Then it can be shown that there exists a Lie algebra sheaf  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  on  $B^{(1)}$ , and a coherent  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)})$  with the action of  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  and canonical isomorphisms  $f^* \mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v) \rightarrow \mathcal{D}_{B'^{(1)}}^1(-\log D'^{(1)} : c_v)$  and  $f^* \mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)}) \rightarrow \mathcal{V}_{\vec{\lambda}}(\mathcal{F}'^{(1)})$ . The objects  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  and  $\mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)})$  do not depend on the choice of  $\mathcal{F}'$ , and these are the objects we wanted to define.

## §6. Local Freeness and Factorization

Here we study the behavior of the  $\mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  module  $\mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)})$  near the discriminant locus  $D^{(1)}$ . Since the problem is local, we can take  $\mathcal{F} = (\pi : C \rightarrow B; s_1, \dots, s_N)$  with condition (Q), with coordinate  $(\tau_1, \dots, \tau_M)$  on  $B$ ,  $M = 3g - 3 + N$ , so that the discriminant locus is of the form  $D = D_1 \cup \dots \cup D_k$ ,  $D_i = \{(\tau) \in B; \tau_i = 0\}, i = 1, \dots, k$ . Set  $E = \bigcap_{j=1}^k D_j$ ,  $E^{(1)} = \bigcap_{j=1}^k D_j^{(1)}$ , and denote by  $\pi_E : C_E \rightarrow E$  the restriction of  $\pi : C \rightarrow B$  on  $E$ . Let  $\tilde{\pi}_E : \tilde{C}_E \rightarrow E$  be the simultaneous normalization of  $\pi_E : C_E \rightarrow E$ , and  $\sigma'_p, \sigma''_p : E \rightarrow \tilde{C}_E$  be the cross-sections corresponding to the normalized double points,  $p = 1, \dots, k$ . Then the family  $\widetilde{\mathcal{F}}_E = \tilde{\pi} : \tilde{C}_E \rightarrow E; \sigma'_p, \sigma''_p, (p = 1, \dots, k), s_1, \dots, s_N$  is a local universal family of  $(N+2k)$ -pointed (not necessary connected fiber) nonsingular curves satisfying the condition (Q). Put  $\widetilde{\mathcal{F}}_E^{(1)} = \widetilde{\pi}_E^{(1)} : \widetilde{C}_E^{(1)} \rightarrow E^{(1)}$ , the associated family of 1-structure. Then the canonical map  $\widetilde{E}^{(1)} \rightarrow E^{(1)}$  is a  $(\mathbf{C}^*)^{2k}$ -principal fibering.

Take  $\omega \in H^0(C \times_B C, \omega_{C/B}^{\boxtimes 2}(2A))$  with the condition (5.7) as well as the property  $i_E^*(\omega) \in H^0(\widetilde{C}_E \times_E \widetilde{C}_E, \omega_{\widetilde{C}_E/E}^{\boxtimes 2}(2A))$ . Using this element  $\omega$ , we fix the trivializations,  $\Theta_{B^{(1)}}(-\log D^{(1)}) \oplus \mathcal{O}_{B^{(1)}} \rightarrow \mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  and  $\Theta_{\widetilde{E}^{(1)}} \oplus \mathcal{O}_{\widetilde{E}^{(1)}} \rightarrow D_{\widetilde{E}^{(1)}}^1(c_v)$ .

Now we have the following fundamental theorem, which we call the sewing process.

**Theorem III.** Fix  $\vec{\lambda} = (\lambda_1, \dots, \lambda_N)$ . Then for each  $\vec{\mu} = (\mu_1, \dots, \mu_k) \in P_\ell^k$ , and for each

$$\Phi \in \text{Hom}_{\mathcal{O}_{\widetilde{E}^{(1)}}(c_v)}(\mathcal{V}_{(\vec{\mu}, \vec{\lambda})}(\widetilde{\mathcal{F}}_E^{(1)}), \mathcal{O}_{\widetilde{E}^{(1)}}) \quad (6.1)$$

we can associate the formal solution having the initial term  $\Phi$ ,

$$\widetilde{\Phi} \in \text{Hom}_{\mathcal{O}_{B^{(1)}}}(\mathcal{V}_{\vec{\lambda}}(\mathcal{F}^{(1)}), \mathcal{O}_{E^{(1)}}[[\tau_1, \dots, \tau_k]]\tau^A \vec{\mu}) \quad (6.2)$$

where  $\tau^{\mu} = \tau_1^{\mu_1}, \dots, \tau_k^{\mu_k}$ ,  $\mu^\dagger = (\mu_1^\dagger, \dots, \mu_k^\dagger)$  and  $\mathcal{O}_{B^{(1)}}[[\tau_1, \dots, \tau_k]] = \varprojlim_m \mathcal{O}_{B^{(1)}} / I_{E^{(1)}}^{m+1}$ .

And we can show that the formal power series  $\tilde{\Phi}$  converge with respect to the variables  $\tau_1, \dots, \tau_k$ .

Then we have

**Theorem IV.** *The coherent  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{V}_\lambda(\mathcal{F}^{(1)})$  is locally free.*

For each  $\mathcal{O}_{B^{(1)}}$ -module  $\mathcal{N}$  such as  $\mathcal{O}_{B^{(1)}}, \Theta_{B^{(1)}}(-\log D^{(1)}), \mathcal{V}_\lambda(\mathcal{F}^{(1)})$ , the  $V$ -filtration along  $E^{(1)}$  is defined as  $V_p \mathcal{N} = I_1^{-p_1} \cdots I_k^{-p_k} \mathcal{N}$ , for  $p = (p_1, \dots, p_k) \in \mathbf{Z}^k$  where  $I_j = \mathcal{O}_{B^{(1)}} \tau_j$ ,  $j = 1, \dots, k$ .

The associated graded module is defined by  $\text{Gr}_*^V \mathcal{N} = \sum_{p \in \mathbf{Z}^k} \text{Gr}_p^V \mathcal{N}$ ,  $\text{Gr}_p^V \mathcal{N} = V_p \mathcal{N} / \sum_{j=1}^k V_{p-e_j} \mathcal{N}$  where  $e_j = (0, \dots, 1, 0, \dots, 0) \in \mathbf{Z}^k$ . Then we have  $\text{Gr}_*^V \mathcal{O}_{B^{(1)}} = \mathcal{O}_{E^{(1)}}[\tau_1, \dots, \tau_k]$ ,  $\text{Gr}_*^V \Theta_{B^{(1)}}(-\log D^{(1)}) = (\sum_{p=1}^k \mathcal{O}_{E^{(1)}} \tau_p \partial / \partial \tau_p + \mathcal{O}_{E^{(1)}}) \otimes \mathbf{C}[\tau_1, \dots, \tau_k]$  where  $\deg \tau_j = -e_j$ , and  $\deg \tau_j \partial / \partial \tau_j = 0$ ,  $j = 1, \dots, k$ .

**Theorem V. 1)** *There exists a canonical  $\mathcal{O}_{\widetilde{E}^{(1)}}$ -module isomorphism*

$$\text{Gr}_0^V \mathcal{V}_\lambda(\mathcal{F}^{(1)}) \otimes_{\mathcal{O}_{E^{(1)}}} \mathcal{O}_{\widetilde{E}^{(1)}} \simeq \sum_{\tilde{\mu} \in P_\ell^k} \mathcal{V}_{(\tilde{\mu}, \tilde{\mu}^\dagger, \lambda)}(\widetilde{\mathcal{F}}_E^{(1)}).$$

2) *There exists a canonical isomorphism of  $\text{Gr}_*^V \Theta_{B^{(1)}}(-\log D^{(1)}) \oplus \text{Gr}_*^V \mathcal{O}_{B^{(1)}}$  modules*

$$\text{Gr}_0^V \mathcal{V}_\lambda(\mathcal{F}^{(1)}) \otimes \mathbf{C}[\tau_1, \dots, \tau_k] \simeq \text{Gr}_*^V \mathcal{V}_\lambda(\mathcal{F}^{(1)}).$$

Let  $k_\theta = \mathbf{C}E_\theta \oplus \mathbf{C}E_{-\theta} \oplus \mathbf{C}H_\theta$  denote the principal 3-dimensional subalgebra of  $\mathfrak{g}$ , and for each  $\lambda \in P_+$ , let  $V_\lambda = \sum_{j \in \frac{1}{2}\mathbf{Z}_{\geq 0}} V_{\lambda, j}$  denote the homogeneous decomposition of the  $k_\theta$  module  $V_\lambda$  into  $(2j+1)$ -dimensional irreducible components.

For each  $(\lambda, \mu, \nu) \in P_\ell^3$ , define

$$\begin{aligned} W_{\lambda, \mu, \nu} &= \{ \varphi \in \text{Hom}_{\mathfrak{g}}(V_\lambda \otimes V_\mu \otimes V_\nu, \mathbf{C}), \\ &\quad \varphi|_{V_{\lambda, j_1} \otimes V_{\mu, j_2} \otimes V_{\nu, j_3}} = 0 \quad \text{if} \quad j_1 + j_2 + j_3 > \ell \} \end{aligned} \tag{6.3}$$

and set  $N_{\lambda, \mu, \nu} = \dim_{\mathbf{C}} W_{\lambda, \mu, \nu}$ . Then we have

**Proposition 6.1.** *For each  $\lambda \in P_\ell^N$ , the rank of  $\mathcal{V}_\lambda(\mathcal{F}^{(1)})$  is computed combinatorially in terms of  $N_{\lambda, \mu, \nu}$ ,  $(\lambda, \mu, \nu) \in P_\ell^3$  only.*

## §7. Monodromy Representations of the Mapping Class Group

Let  $\Gamma_{g,N}$  be the mapping class group of  $N$ -pointed oriented surfaces of genus  $g$  with first-order infinitesimal structure, and put  $\Gamma_g = \Gamma_{g,0}$ . Then there exist surjective homomorphisms,  $\Gamma_{g,N} \rightarrow S_N \rightarrow 1$  and  $\Gamma_{g,N} \rightarrow \Gamma_g \rightarrow 1$ , and also an

exact sequence  $1 \rightarrow \mathbf{Z}^N \rightarrow \Gamma_{g,N} \rightarrow \Gamma_g^N \rightarrow 1$ , where  $\Gamma_g^N$  is the mapping class group of  $N$ -pointed oriented surfaces of genus  $g$ .

Next consider the modular stack  $M_{g,N}^{(1)}$  of the  $N$ -pointed (unordered) stable curves of genus  $g$  with first-order infinitesimal structure, and the normal crossing divisor  $D_{g,N}^{(1)} \subset M_{g,N}^{(1)}$  representing singular curves. Set  $\overset{\circ}{M}_{g,N}^{(1)} = M_{g,N}^{(1)} - D_{g,N}^{(1)}$ . Then the fundamental group  $\pi_1(\overset{\circ}{M}_{g,N}^{(1)})$  is isomorphic to  $\Gamma_{g,N}$ . For an  $S_N$ -invariant subset  $A$  of  $P_\ell^N$ , we denote  $\mathcal{V}_A(\mathcal{F}^{(1)}) = \bigoplus_{\lambda \in A} \mathcal{V}_\lambda(\mathcal{F}^{(1)})$ . Then the functors  $\mathcal{F} \rightarrow \mathcal{D}_{B^{(1)}}^1(-\log D^{(1)} : c_v)$  and  $\mathcal{F} \rightarrow \mathcal{V}_A(\mathcal{F}^{(1)})$  define a sheaf of twisted first-order differential operators  $\mathcal{D}_{M_{g,N}^{(1)}}^1(-\log D_{g,N}^{(1)} : c_v)$  on  $M_{g,N}^{(1)}$ , and a locally-free coherent  $\mathcal{O}_{M_{g,N}^{(1)}}$ -module  $\mathcal{V}_A(M_{g,N}^{(1)})$  on which  $\mathcal{D}_{M_{g,N}^{(1)}}^1(-\log D_{g,N}^{(1)} : c_v)$  acts. The functor  $(\mathcal{F}, \omega) \rightarrow \mathcal{L}_\omega(\mathcal{F}^{(1)})$  define a system of invertible  $\mathcal{O}_{U_i}$ -modules  $\mathcal{L}_i$  on which  $\mathcal{D}_{U_i}^1(-\log D^{(1)} : c_v)$  acts for some open covering  $\{U_i\}$  of  $M_{g,N}^{(1)}$ . Now restrict the system  $\{\mathcal{L}_i, U_i\}$  on  $\overset{\circ}{M}_{g,N}^{(1)}$ . Then the twisting system  $\{\mathcal{L}_{ij}, U_i \cap U_j\}$  on  $\overset{\circ}{M}_{g,N}^{(1)}$  is defined by  $\mathcal{L}_{ij} = \mathcal{H}\text{om}_{\mathcal{D}_{M_{g,N}^{(1)}}^1(c_v)}(\mathcal{L}_{i|U_i \cap U_j}, \mathcal{L}_{j|U_i \cap U_j})$ , cf. Kashiwara [7]. Then using some results of Oda [10], about the homotopy type of  $\overset{\circ}{M}_{g,N}^{(1)}$ , we have

**Proposition 7.1.** *The twisting data  $\{\mathcal{L}_{ij}\}$  on  $\overset{\circ}{M}_{g,N}^{(1)}$  define central extensions  $\widehat{\Gamma}_{g,N}$  and  $\widehat{\Gamma}_g$  of  $\Gamma_{g,N}$  and  $\Gamma_g$  by the multiplicative subgroup  $K$  of  $\mathbf{C}^*$  generated by  $e^{2\pi\sqrt{-1}c_v/24}$ , with the following commutative diagram.*

$$\begin{array}{ccccccc} 1 & \longrightarrow & K & \longrightarrow & \widehat{M}_{g,N} & \longrightarrow & 1 \\ & & \downarrow = & & \downarrow & & \\ 1 & \longrightarrow & K & \longrightarrow & \widehat{M}_g & \longrightarrow & 1. \end{array} \quad (7.1)$$

**Theorem VI.** *For each  $A$ , the system of solution sheaves on  $\overset{\circ}{M}_{g,N}^{(1)}$ ,*

$$\mathcal{H}\text{om}_{\mathcal{D}_{M_{g,N}^{(1)}}^1(c_v)}(\mathcal{V}_A|_{U_i}, \mathcal{L}_i) \quad (7.2)$$

*defines a monodromy representation  $\varrho$  of the group  $\widehat{\Gamma}_{g,N}$  such that an element  $k$  of  $K$  acts as  $\varrho(k) = k \cdot \text{id}$ .*

## References

1. Belavin, A.A., Polyakov, A.M., Zamolodchikov, A.B.: Infinite conformal symmetry in two dimensional quantum field theory. *Nucl. Phys.* **B241** (1984) 333–380
2. Beilinson, A.A., Schechtman, Y.A.: Determinant bundles and Virasoro algebras. *Commun. Math. Phys.* **118** (1988) 651–701
3. Deligne, P., Mumford, D.: The irreducibility of the space of curves of given genus. *Publ. Math. I.H.E.S.* **36** (1969) 75–110

4. Friedan, D., Shenker, S.: The analytic geometry of two dimensional conformal field theory. *Nucl. Phys.* **B281** (1987) 509–545
5. Gabber, O.: The integrability of the characteristic variety. *Amer. J. Math.* **103** (1981) 445–468
6. Kac, V.G.: Infinite dimensional Lie algebras. 2nd ed. Cambridge University Press, 1985
7. Kashiwara, M.: Representation theory and  $\mathcal{D}$ -modules on flag varieties. *Astérisque* **173-174** (1989) 55–109
8. Knizhnik, V.G., Zamolodchikov, A.B.: Current algebra and Wess-Zumino model in two dimensions. *Nucl. Phys.* **B247** (1984) 83–103
9. Moore, G., Seiberg, N.: Classical and quantum conformal field theory. *Commun. Math. Phys.* **123** (1989) 177–254
10. Oda, T.: Etale homotopy type of the moduli space of algebraic curves. Preprint, R.I.M.S., 1990
11. Segal, G.: Conformal field theory. In these Proceedings, pp. 1387–1396
12. Tsuchiya, A., Kaneko, Y.: Vertex operators in conformal field theory on  $P^1$  and monodromy representations of the braid group. In: Conformal field theory and solvable lattice models. *Adv. Stud. Pure Math.* **16** (1988) 297–372
13. Tsuchiya, A., Ueno, K., Yamada, Y.: Conformal field theory on universal family of stable curves with gauge symmetries. In: Integrable systems in quantum field theory and statistical mechanics. *Adv. Stud. Pure Math.* **19** (1989) 459–566
14. Verlinde, E.: Fusion rules and modular transformations in 2d conformal field theory. *Nuclear Phys.* **B300** (1988) 360–376
15. Witten, E.: Quantum field theory and the Jones polynomial. *Commun. Math. Phys.* **121** (1989) 351–399



# Non-Constructive Proofs in Combinatorics

Noga Alon

Department of Mathematics, Raymond and Beverly Sackler Faculty of Exact Sciences  
Tel Aviv University, Tel Aviv, Israel and  
IBM Almaden Research Center, San Jose, CA 95120, USA

One of the main reasons for the fast development of Combinatorics during the recent years is certainly the widely used application of combinatorial methods in the study and the development of efficient algorithms. It is therefore somewhat surprising that many results proved by applying some of the modern combinatorial techniques, including Topological methods, Algebraic methods, and Probabilistic methods, merely supply existence proofs and do not yield efficient (deterministic or randomized) algorithms for the corresponding problems.

We describe some representing non-constructive proofs of this type, demonstrating the applications of Topological, Algebraic and Probabilistic methods in Combinatorics, and discuss the related algorithmic problems.

## 1. Topological Methods

The application of topological methods in the study of combinatorial objects like partially ordered sets, graphs, hypergraphs and their coloring have become in the last ten years part of the mathematical machinery commonly used in combinatorics. Many interesting examples appear in [12]. Some of the more recent results of this type deal with problems that are closely related to certain algorithmic problems. While the topological tools provide a powerful technique for proving the required results, they give us no clue on an efficient way for solving the corresponding algorithmic questions.

A typical result of this type is the following theorem, proved in [2].

**Theorem 1.1.** *Let  $N$  be an open necklace with  $ka_i$  beads of color  $i$ ,  $1 \leq i \leq t$ . Then one can cut  $N$  in  $(k - 1)t$  places and partition the resulting intervals into  $k$  collections, each containing precisely  $a_i$  beads of color  $i$  for all  $1 \leq i \leq t$ .*

The bound  $(k - 1)t$ , conjectured in [17] (where it is proved for  $k = 2$ ) is sharp. This can be seen by considering the necklace in which the beads of each type appear contiguously. The proof supplies no efficient procedure, which finds, given a necklace as above, a partition of it with the desired properties. By an efficient procedure

we mean here, and in what follows, either a deterministic algorithm whose running time is polynomial (in the length of the input) or a randomized algorithm whose expected running time (on the worst-case input) is polynomial.

Here is a sketch of the proof of the above theorem. A similar method is used in [6]. First we need a continuous version of it. Let  $I = [0, 1]$  be the (closed) unit interval. An *interval  $t$ -coloring* is a coloring of the points of  $I$  by  $t$  colors, such that for each  $i$ ,  $1 \leq i \leq t$ , the set of points colored  $i$  is (Lebesgue) measurable. Given such a coloring, a  $k$ -*splitting of size  $r$*  is a sequence of numbers  $0 = y_0 \leq y_1 \leq \dots \leq y_r \leq y_{r+1} = 1$  and a partition of the family of  $r + 1$  intervals  $F = \{[y_i, y_{i+1}] : 0 \leq i \leq r\}$  into  $k$  pairwise disjoint subfamilies  $F_1, \dots, F_k$  whose union is  $F$ , such that for each  $j$ ,  $1 \leq j \leq k$ , the union of all intervals in  $F_j$  captures precisely  $1/k$  of the total measure of each of the  $t$  colors.

The following result is the continuous analogue of Theorem 1.1.

**Proposition 1.2.** *Every interval  $t$ -coloring has a  $k$ -splitting of size  $(k - 1)t$ .*

We note that a similar statement can be proved for general continuous probability measures instead of those defined by the colors. This generalizes the Hobby-Rice Theorem on  $L_1$ -approximation [18]. It is also related to one of the cake-splitting problems of Steinhaus. It is easy to see that the classical theorem of Liapounoff [20] implies the existence of an even splitting in this more general setting, but unlike the above result does not supply any finite bound on the number of cuts required to form the splitting. For more details see [2].

It is not difficult to see that Proposition 1.2 implies Theorem 1.1. This is because any open necklace with  $\sum_{i=1}^t ka_i = kn$  beads as in the theorem can be converted into an interval  $t$ -coloring by partitioning the interval  $I$  into  $kn$  segments of equal size and by coloring the  $j$ -th part by the color of the  $j$ -th bead of the necklace. By Proposition 1.2 there is a  $k$  splitting with  $(k - 1)t$  cuts. Of course, these cuts need not occur at the endpoints of the segments, but a simple induction argument can be used to show that the cuts may be shifted until they form a partition of the discrete necklace satisfying the assertion of Theorem 1.1. We omit the details.

Another simple observation, whose details we omit, is the fact that the validity of Proposition 1.2 for  $(t, k)$  and for  $(t, k')$  implies its validity for  $(t, kk')$ . Therefore, it suffices to prove the proposition for prime values of  $k$ . To do so we define, following [11], a CW-complex  $Y = Y(k, m)$  as follows.

For two integers  $k$  and  $m$ , put  $N = N(k, m) = (k - 1)(m + 1)$  and let  $\Delta = \Delta^N$  denote the  $N$ -dimensional simplex; i.e.,  $\Delta = \{(x_0, \dots, x_N) : x_i \geq 0, \sum_{i=0}^N x_i = 1\}$ . The *support* of a point  $x \in \Delta$ , denoted by  $\text{Supp}(x)$ , is the minimal face of  $\Delta$  that contains  $x$ . Define

$$\begin{aligned} Y = Y(k, m) &= \{(y_1, \dots, y_k) : y_1, \dots, y_k \in \Delta, \text{Supp}(y_i) \cap \text{Supp}(y_j) \\ &= \emptyset \text{ for all } 1 \leq i < j \leq k\}. \end{aligned}$$

The cyclic group  $Z_k$  acts freely on  $Y$  by letting its generator  $\omega$  cyclically shift the coordinates of each point  $y \in Y$ , i.e.,  $\omega(y_1, \dots, y_k) = (y_2, \dots, y_k, y_1)$ .

The following lemma is proved in [11].

**Lemma 1.3.** *If  $k$  is a prime,  $m \geq 1$ ,  $N = N(k, m) = (k - 1)(m + 1)$  and  $Y = Y(k, m)$  and  $\omega$  are as in the preceding paragraph, then  $Y$  is  $N - k$  connected and hence for every continuous mapping  $h : Y \mapsto R^m$  there is a whole orbit of the  $Z_k$  action on  $Y$  that is mapped by  $h$  into one point. I.e., there is a  $y \in Y$  such that  $h(y) = h(\omega(y)) = \dots = h(\omega^{k-1}(y))$ .*

We can now prove Proposition 1.2 for primes  $k$ . Let  $c$  be an interval  $t$ -coloring. Define  $N = N(k, t - 1) = (k - 1)t$ ,  $Y = Y(k, t - 1)$  and consider the continuous function  $h : Y \mapsto R^{t-1}$  defined as follows.

Suppose  $y = (y_1, \dots, y_k) \in Y$ . By the definition of  $Y$ , each  $y_i$  is a point of  $A^N$ , i.e., a real vector of length  $N$  with nonnegative coordinates whose sum is 1. Moreover, the supports of the points  $y_i$  are pairwise disjoint. Put  $x = (x_0, \dots, x_N) = \frac{1}{k} \sum_{i=1}^k y_i$ , and define a partition of the  $[0, 1]$ -interval  $I$  into  $N + 1$  intervals  $I_0, \dots, I_N$  by

$$I_0 = [0, x_0] \quad \text{and} \quad I_j = \left[ \sum_{l=0}^{j-1} x_l, \sum_{l=0}^j x_l \right], \quad (1 \leq j \leq N).$$

Observe that since the supports of the points  $y_i$  are pairwise disjoint, then for each interval  $I_j$  with a positive length there is a unique  $l$  such that the  $j$ -th coordinate of  $y_l$  is positive.

For each  $l$ ,  $1 \leq l \leq k$ , let  $F_l$  be the family of all the intervals  $I_j$  such that the  $j$ -th coordinate of  $y_l$  is positive. Note that the sum of lengths of the intervals in each  $F_l$  is precisely  $1/k$ , and that  $F_1, \dots, F_k$  form a partition of all the intervals  $I_j$  whose lengths are positive. For each  $i$ ,  $1 \leq i \leq t - 1$ , define  $h_i(y)$  to be the measure of the  $i$ -th color in the union of the intervals of  $F_1$ . The function  $h(y)$  is now defined by  $h(y) = (h_1(y), \dots, h_{t-1}(y))$ .

This function is clearly continuous. Also, for every  $1 \leq l \leq k$  and  $1 \leq i \leq t - 1$ ,  $h_i(\omega^{l-1}(y))$  is precisely the measure of the  $i$ -th color in the union of the intervals of  $F_l$ . By Lemma 1.3 there is a  $y \in Y$  such that  $h(y) = h(\omega(y)) = \dots = h(\omega^{k-1}(y))$ . This means that each of the  $k$  families  $F_l$  corresponding to this point  $y$  captures precisely  $1/k$  of the measure of each of the first  $t - 1$  colors. Since the total measure of each  $F_l$  is  $1/k$ , it follows that the last color is evenly distributed between the families as well. This completes the proof for the case of prime  $k$ , and hence implies the validity of Proposition 1.2 and Theorem 1.1.  $\square$

The main topological tool in the above proof is the Borsuk-type theorem stated in Lemma 1.3. This proof does not seem to supply an efficient way of producing a partition whose existence is guaranteed by the theorem.

In the classification of algorithmic problems according to their complexity, it is customary to try and identify the problems that can be solved efficiently, and those that *probably* cannot be solved efficiently. A class of problems that can be solved

efficiently is the class  $P$  of all problems for which there are deterministic algorithms whose running time is polynomial in the length of the input. A class of problems that probably cannot be solved efficiently are all the  $NP$ -complete problems. An extensive list of such problems appears in [16]. It is well known that if any of them can be solved efficiently, then so can all of them, since this would imply that the two complexity classes  $P$  and  $NP$  are equal.

It is not too difficult to show that the following problem is  $NP$ -complete: Given a necklace satisfying the assumptions of Theorem 1.1, decide if one can form an even  $k$ -splitting of it by using less than  $b$  cuts. On the other hand, we know that  $(k - 1)t$  cuts always suffice, so although the problem of finding the minimum possible number of cuts cannot be solved efficiently, unless  $P = NP$ , it is considerable and seems likely that the problem of finding an even  $k$ -splitting using  $(k - 1)t$  cuts is much easier. We do not know any efficient algorithm for this problem.

Another result whose (simple) proof applies the Borsuk-Ulam theorem is the following fact, proved in [1]:

**Theorem 1.4.** *Let  $A_1, \dots, A_d$  be  $d$  pairwise disjoint subsets of  $\mathbb{R}^d$ , each containing precisely  $n$  points, and suppose that no hyperplane contains  $d + 1$  of the points in the union of all the sets  $A_j$ . Then there is a partition of  $\bigcup A_j$  into  $n$  pairwise disjoint sets  $S_1, \dots, S_n$ , each containing precisely one point from each  $A_j$ , such that the  $n$  simplices  $\text{conv}(S_1), \dots, \text{conv}(S_n)$  are pairwise disjoint.*

Here, again, the proof does not supply an efficient way of finding the sets  $S_i$  if the sets  $A_j$  are given, (although the proof does provide an efficient way of doing it for each *fixed* dimension  $d$ .)

## 2. Algebraic Methods

Many combinatorial proofs rely on methods from linear and multilinear algebra. Extensive survey of results of this type is given in [9]. These proofs rarely supply constructive procedures for the corresponding algorithmic problems. Here is a simple example, which is a special case of one of the results in [5].

**Proposition 2.1.** *Every (not necessarily simple) graph with maximum degree 5 and average degree greater than 4, contains a 3-regular subgraph.*

The proof relies on the classical theorem of Chevalley and Warning (see, e.g., [10]). This theorem, that deals with the number of solutions of a system of multi-variable polynomials over a finite field, is the following.

**Theorem 2.2.** *Let  $P_j(x_1, \dots, x_m)$ ,  $(1 \leq j \leq n)$  be  $n$  polynomials over a finite field  $F$  of characteristic  $p$ . If the number of variables,  $m$ , is greater than the sum of the degrees of the polynomials then the number of common zeros of the polynomials (in  $F^m$ ) is divisible by  $p$ . In particular, if there is one common zero then there is another one.*

The proof is extremely simple; If  $F$  has  $q$  elements, then the number  $N$  of common zeros satisfies

$$N \equiv \sum_{x_1, \dots, x_m \in F} \prod_{j=1}^n (1 - P_j(x_1, \dots, x_m)^{q-1}) \pmod{p}.$$

By expanding the right hand side we get a linear combination of monomials of the form  $\prod_{i=1}^m x_i^{k_i}$  and for each such monomial at least one of the exponents  $k_i$  is strictly smaller than  $q - 1$ . This implies that in  $F$ ,  $\sum_{x_i \in F} x_i^{k_i} = 0$ , showing that the contribution of each monomial to the sum expressing  $N$  is  $0 \pmod{p}$  and completing the proof.  $\square$

We can now prove Proposition 2.1. Given a graph  $G = (V, E)$  satisfying the assumptions of the proposition, let  $n$  denote the number of its vertices. For each edge  $e \in E$  and for each vertex  $v \in V$ , let  $a(v, e)$  be 0 if  $e$  is not incident with  $v$ , 1 if  $e$  is a non-loop incident with  $v$ , and 2 if  $e$  is a loop incident with  $v$ . For each  $e \in E$  let  $x_e$  be a variable and consider the following system of polynomial equations over  $GF(3)$ :

$$\sum_{e \in E} a(v, e) x_e^2 = 0 \quad (v \in V).$$

This is a system of  $n$  degree-2 polynomial equations with  $|E| > 2n$  variables. Moreover, it clearly has the trivial solution  $x_e = 0$  for all  $e$ . Hence there is, by Theorem 2.2, a non-trivial solution  $(y_e : e \in E)$ . Let  $H$  be the subgraph of  $G$  consisting of all edges  $e$  for which  $y_e \neq 0$ . By the equations above, the degree of every vertex of  $H$  is divisible by 3, and since the maximum degree in  $G$  is 5 it follows that  $H$  is 3-regular, completing the proof.  $\square$

It is known that the decision problem: “Given a graph  $G$ , decide if it contains a 3-regular subgraph”, is  $NP$ -complete. By the proposition above in certain cases we know that the answer to the decision problem is “yes” and yet the proof does not yield an efficient procedure for finding such a subgraph.

Another result proved by applying some extension of the Chevalley Warning Theorem is the following statement, proved in [7]. Recall that a *hypergraph* is a pair  $(V, \mathcal{F})$  (sometimes denoted only by  $\mathcal{F}$ ), where  $V$  is a finite set of vertices, and  $\mathcal{F}$  is a finite set of subsets of  $V$ . The *degree* of a vertex is the number of edges that contain it.

**Theorem 2.3.** *Let  $q$  be a prime power, and let  $\mathcal{F} = \{F_1, \dots, F_{d(q-1)+1}\}$  be a hypergraph whose maximal degree is  $d$ . Then there exists  $\emptyset \neq \mathcal{F}_0 \subset \mathcal{F}$ , such that  $|\bigcup_{F \in \mathcal{F}_0} F| \equiv 0 \pmod{q}$ .*

Here, again, we do not know how to quickly find such a subset  $\mathcal{F}_0$ . Moreover, it can be shown that the problem of finding such an  $\mathcal{F}_0$  is equivalent to the following problem: Given a polynomial  $h$  of degree at most  $d$  with  $d(q-1)+1$  variables over  $GF(q)$ , suppose that  $h(0) = 0$ . Find another zero of  $h(x)$  in which each variable is either 0 or 1.

### 3. Probabilistic Methods

Probabilistic methods have been useful in combinatorics for almost fifty years. Many examples can be found in [14] and in [21].

In a typical application of the probabilistic method we try to prove the existence of a combinatorial structure (or a substructure of a given structure) with certain prescribed properties. To do so, we show that a randomly chosen element from an appropriately defined sample space satisfies all the required properties with positive probability. In most applications, this probability is not only positive, but is actually high and frequently tends to 1 as the parameters of the problem tend to  $\infty$ . In such cases, the proof usually supplies an efficient randomized algorithm for producing a structure of the desired type, and in many cases this algorithm can be derandomized and converted into an efficient deterministic one.

There are, however, certain examples, where one can prove the existence of the required combinatorial structure by probabilistic arguments that deal with rare events; events that hold with positive probability which is exponentially small in the size of the input. Such proofs usually yield neither randomized nor deterministic efficient procedures for the corresponding algorithmic problems.

A class of examples demonstrating this phenomenon is the class of results proved by applying the Local Lemma. This result, proved in [13] (see also, e.g., [21]), supplies a way of showing that certain events hold with positive probability, although this probability may be extremely small. The exact statement (for the symmetric case) is the following.

**Lemma 3.1.** *Let  $A_1, \dots, A_n$  be events in an arbitrary probability space. Suppose that the probability of each of the  $n$  events is at most  $p$ , and suppose that each event  $A_i$  is mutually independent of all but at most  $b$  of the other events  $A_j$ . If  $ep(b + 1) < 1$  then with positive probability none of the events  $A_i$  holds.*

One of the applications of this lemma, given already in the original paper [13], deals with *hypergraph coloring*. A hypergraph is *k-uniform* if each of its edges contains precisely  $k$  vertices. It is *k-regular* if each of its vertices is contained in precisely  $k$  edges. A hypergraph is *2-colorable* if there is a two-coloring of the set of its vertices so that none of its edges is monochromatic. Erdős and Lovász proved the following result.

**Proposition 3.2.** *For each  $k \geq 9$ , every  $k$ -regular,  $k$ -uniform hypergraph is two colorable.*

The proof follows almost immediately from Lemma 3.1. Let  $(V, E)$  be a  $k$ -uniform,  $k$ -regular hypergraph, and let  $f: V \mapsto \{0, 1\}$  be a random 2-coloring obtained by choosing, for each  $v \in V$  randomly and independently,  $f(v) \in \{0, 1\}$  according to a uniform distribution. For each  $e \in E$  let  $A_e$  denote the event that  $f$  restricted to  $e$  is a constant, i.e., that  $e$  is monochromatic. It is obvious that  $\text{Prob}(A_e) = 2^{-(k-1)}$  for every  $e$ , and that each event  $A_e$  is mutually independent of all the events  $A_f$  but those for which  $f \cap e \neq \emptyset$ . Since there are at most  $k(k - 1)$

edges  $f$  that intersect  $e$  we can substitute  $b = k(k - 1)$  and  $p = 2^{-(k-1)}$  in Lemma 3.1 and conclude that for  $k \geq 9$  with positive probability none of the events  $A_e$  holds, completing the proof.  $\square$

We note that a different, algebraic proof of the statement of the last proposition (that works for all  $k \geq 8$ ) is given in [4]. Both proofs do not supply an efficient way of finding a proper two-coloring for a given hypergraph satisfying the assumptions of the proposition. Note that, in general, the problem of deciding whether a hypergraph is 2-colorable is *NP*-complete.

Another application of the Local Lemma, which appears in [8], is the following.

**Proposition 3.3.** *Every directed simple graph  $D = (V, E)$  with minimum outdegree  $\delta$  and maximum indegree  $A$  contains a directed (simple) cycle of length  $0(\bmod k)$ , provided  $e(A\delta + 1)\left(1 - \frac{1}{k}\right)^{\delta} < 1$ .*

The proof here first applies the Local Lemma to show that there exists a function  $f: V \mapsto \{0, 1, \dots, k\}$  such that for every  $v \in V$  there is a vertex  $u \in V$  such that  $(v, u)$  is a directed edge of  $D$  and  $f(u) \equiv f(v) + 1(\bmod k)$ .

Given such an  $f$ , the rest of the proof is very simple. We just choose, for every vertex  $v \in V$ , some vertex  $p(v)$  such that  $(v, p(v))$  is a directed edge and  $f(p(v)) \equiv f(v) + 1(\bmod k)$ . Suppose  $v \in V$  and consider the sequence

$$v_0 = v, \quad v_1 = p(v_0), \quad v_2 = p(v_1), \dots$$

Let  $j$  be the minimum index such that there is an  $i < j$  with  $v_i = v_j$ . The cycle  $v_iv_{i+1} \dots v_{j-1}v_j = v_i$  is a directed cycle of length  $0(\bmod k)$ , as needed.

Here, again, the proof is not constructive in the sense that it does not provide an efficient way of finding such a cycle in a directed graph satisfying the assumptions. This is because the proof that a function  $f$  as above exists is non-constructive.

We note that it is not known if the related decision problem “Given a directed graph, decide if it contains a directed even cycle” is polynomial, but it is easy to deduce from the results of [15] that the similar problem “Given a directed graph and an edge  $e$  in it, decide if there is an even cycle containing  $e$ ” is *NP*-complete.

The proof of the next result also relies on the Local Lemma, but contains several additional ingredients as well. The details appear in [3].

**Theorem 3.4.** *There is an absolute constant  $c$  with the following property: For any two graphs  $G_1 = (V, E_1)$  and  $G_2 = (V, E_2)$  on the same set of vertices, where  $G_1$  has maximum degree at most  $d$  and  $G_2$  is a vertex disjoint union of cliques of size  $cd$  each, the chromatic number of the graph  $G = (V, E_1 \cup E_2)$  is precisely  $cd$ .*

The proof, again, does not supply an efficient (deterministic or randomized) algorithm for producing a proper  $cd$ -vertex coloring of  $G$ .

We close this section mentioning the following result of J. Spencer, whose proof, given in [22], which combines the probabilistic method with a counting argu-

ment, also fails to supply an efficient procedure for the corresponding algorithmic problem.

**Theorem 3.5.** *Let  $v_1, \dots, v_n$  be  $n$  real vectors of length  $n$  each, and suppose that the  $l_\infty$ -norm of each  $v_i$  is at most 1. Then there are  $\varepsilon_1, \dots, \varepsilon_n \in \{-1, 1\}$ , such that the  $l_\infty$ -norm of the sum  $\sum_{i=1}^n \varepsilon_i v_i$  is at most  $6\sqrt{n}$ .*

## 4. Concluding Remarks

We have seen several examples of combinatorial results proved by topological, algebraic or probabilistic methods. One natural question that arises is whether these methods are necessary. After all, we may tend to believe that simply stated combinatorial results should have simple combinatorial proofs. Although this sounds plausible, there are no known natural combinatorial proofs for any of the results mentioned here (as well as for various other known similar examples).

Another question that should be addressed is whether the proofs given here are really inherently non-constructive. Is it possible to modify them so that they yield efficient ways of solving the corresponding algorithmic problems? There are no known efficient algorithms for any of the problems mentioned here. However, it seems very likely that such algorithms do exist. This is related to questions regarding the complexity of search problems that have been studied by several researchers. See, e.g., [19].

In the study of complexity classes like  $P$  and  $NP$  one usually considers only decision problems, i.e., problems for which the only two possible answers are “yes” or “no.” However, the definitions extend easily to the so called “search” problems, which are problems where a more elaborate output is sought. The search problems corresponding to the complexity classes  $P$  and  $NP$  are sometimes denoted by  $FP$  and  $FNP$ .

Consider, for example, the obvious algorithmic problem suggested by Theorem 1.1, namely, given a necklace satisfying the assumptions of the theorem, find a partition of it satisfying the conclusions of the theorem. This problem is in  $FNP$ , since it is a search problem, and given a proposed solution for it we can check in polynomial time that it is indeed a solution.

Notice that this problem always has a solution, by Theorem 1.1, and hence it seems plausible that finding one should not be a very difficult task. The situation is similar with all the other algorithmic problems corresponding to the various results mentioned here. Still, the problem of solving efficiently the corresponding search problems remains an intriguing open question.

## References

1. J. Akiyama, N. Alon: Disjoint simplices and geometric hypergraphs. In: Combinatorial Mathematics; Proc. of the Third International Conference, New York, NY 1985 (G.S. Blum, R.L. Graham and J. Malkevitch, eds.). Ann. New York Acad. Sci. **555** (1989) 1–3
2. N. Alon: Splitting necklaces. Adv. Math. **63** (1987) 247–253

3. N. Alon: The strong chromatic number of a graph. To appear
4. N. Alon, Z. Bregman: Every 8-uniform, 8-regular hypergraph is 2-colorable. *Graphs and Combinatorics* **4** (1988) 303–305
5. N. Alon, S. Friedland, G. Kalai: Regular subgraphs of almost regular graphs. *J. Combin. Theory Ser. B* **37** (1984) 79–91
6. N. Alon, P. Frankl, L. Lovász: The chromatic number of Kneser hypergraphs. *Trans. Amer. Math. Soc.* **298** (1986) 359–370
7. N. Alon, D.J. Kleitman, R. Lipton, R. Meshulam, M.O. Rabin, J. Spencer: Set systems with no union of cardinality 0 modulo  $m$ . To appear
8. N. Alon, N. Linial: Cycles of length 0 modulo  $k$  in directed graphs. *J. Combin. Theory Ser. B* **47** (1989) 114–119
9. L. Babai, P. Frankl: Linear algebra methods in combinatorics I. Preliminary version. Department of Computer Science, University of Chicago, 1988
10. Z.I. Borevich, I.R. Shafarevich: Number theory. Academic Press, New York 1966
11. I. Bárány, S.B. Shlosman, A. Szűcs: On a topological generalization of a theorem of Tverberg. *J. London Math. Soc.* (2), **23** (1981) 158–164
12. A. Björner: Topological methods. To appear in: *Handbook of Combinatorics* (R.L. Graham, M. Grötschel and L. Lovász, eds.). North-Holland, Amsterdam 1991
13. P. Erdős, L. Lovász: Problems and results on 3-chromatic hypergraphs and some related questions. In: *Infinite and finite sets* (A. Hajnal et. al., eds.). Colloq. Math. Soc. J. Bolyai, vol. 11. North-Holland, Amsterdam 1975, pp. 609–627
14. P. Erdős, J. Spencer: Probabilistic methods in combinatorics. Academic Press/Akadémiai Kiadó, New York/Budapest 1974
15. S. Fortune, J.E. Hopcroft, J. Wyllie: The directed subgraph homeomorphism problem. *Theor. Comp. Sci.* **10** (1980) 111–120
16. M.R. Garey, D.S. Johnson: Computers and intractability: A guide to the theory of NP-completeness. Freeman, 1979
17. C.H. Goldberg, D.B. West: Bisection of circle colorings. *SIAM J. Algeb. Discrete Methods* **6** (1985) 93–106
18. C.R. Hobby, J.R. Rice: A moment problem in  $L_1$  approximation. *Proc. Amer. Math. Soc.* **16** (1965) 665–670
19. D.S. Johnson, C.H. Papadimitriou, M. Yannakakis: How easy is local search? *JCSS* **37** (1988) 79–100
20. A. Liapounoff: Sur les fonctions vecteurs complètement additives. *Izv. Akad. Nauk SSSR* **4** (1940) 465–478
21. J. Spencer: Ten lectures on the probabilistic method. SIAM, Philadelphia, 1987
22. J. Spencer: Six standard deviations suffice. *Trans. Amer. Math. Soc.* **289** (1985) 679–706



# Infinite Permutation Groups in Enumeration and Model Theory

Peter J. Cameron

School of Mathematical Sciences, QMW, University of London, Mile End Road  
London E1 4NS, UK

## 1. Introduction

A permutation group  $G$  on a set  $\Omega$  has a natural action on  $\Omega^n$  for each natural number  $n$ . The group is called *oligomorphic* if it has only finitely many orbits on  $\Omega^n$  for all  $n \in \mathbb{N}$ . (The term means “few shapes”. Typically our permutation groups are groups of automorphisms of structures of some kind; oligomorphy implies that the structure has only finitely many non-isomorphic  $n$ -element substructures for each  $n$ .)

Oligomorphic permutation groups have close connections with both model theory and combinatorial enumeration. For the former, a basic result is the theorem of Engeler, Ryll-Nardzewski and Svenonius characterizing  $\aleph_0$ -categorical countable structures as those whose automorphism groups are oligomorphic. The connection with enumeration is via homogeneous structures, those for which orbits on  $n$ -sets are isomorphism types of induced substructures. A theorem of Fraïssé gives us a rich supply of homogeneous structures. These matters are described in Section 2. Section 3 develops some tools of enumeration theory (cycle index) in this context, with a few applications. In Section 4, the famous countable “random graph” of Erdős and Rényi is used to introduce the ideas of measure and Baire category. In the fifth section, some results and problems on the rate of growth of orbit-counting sequences are presented. Finally, the search for cyclic automorphisms of certain interesting graphs leads to sum-free sets, which have a fascinating theory.

For an oligomorphic permutation group  $G$ , I let  $f_n(G)$ ,  $F_n(G)$  and  $F_n^*(G)$  be the numbers of orbits of  $G$  on  $n$ -sets,  $n$ -tuples of distinct elements, and all  $n$ -tuples respectively. By convention,  $f_0(G) = F_0(G) = F_0^*(G) = 1$ . There are some interesting relations among these sequences, notably

$$F_n^*(G) = \sum_{k=1}^n S(n, k) F_k(G),$$

where  $S(n, k)$  is the Stirling number of the second kind; this fact has a number of combinatorial consequences (Cameron and Taylor 1985).

I will always assume that the set  $\Omega$  on which a group acts is finite or countable. (Little is lost here; the downward Löwenheim–Skolem theorem of model theory guarantees that any sequences  $(f_n(G))$ ,  $(F_n(G))$ , etc. realized by an oligomorphic group can be realized by a group of countable degree.)

*Notation:*  $G \times H$  and  $G \text{ Wr } H$  denote the direct and wreath products of the permutation groups  $G$  and  $H$  (acting on the disjoint union and cartesian product respectively of the sets admitting  $G$  and  $H$ );  $G_\alpha$  is the stabilizer of the point  $\alpha \in \Omega$ .

$S_\infty$  is the symmetric group of countable degree;  $C_n$ , the cyclic group of order  $n$ ; and  $A$ , the group of order-preserving permutations of  $\mathbb{Q}$ . Note that  $f_n(A) = 1$  for all  $n$ : any order-preserving bijection between finite subsets of  $\mathbb{Q}$  can be extended to a (piecewise-linear) order-preserving permutation of  $\mathbb{Q}$ .

Three properties of  $\mathbb{Q}$  will be important. First of course is Cantor's (1895) characterization of  $\mathbb{Q}$  as countable dense ordered set without endpoints. Second, as just mentioned, any order-preserving bijection between finite subsets of  $\mathbb{Q}$  extends to an automorphism. Third is the fact that, if  $X, Y$  are finite ordered sets with  $X \subseteq Y$ , then any embedding of  $X$  in  $\mathbb{Q}$  can be extended to an embedding of  $Y$ . These observations are the starting point for the next section.

For a fuller and more leisurely discussion of oligomorphic permutation groups, see my lecture notes (Cameron 1990).

## 2. $\aleph_0$ -Categoricity and Homogeneity

A countable structure  $M$  over a first-order language is  $\aleph_0$ -categorical if it is the unique countable model of its theory, i.e. determined up to isomorphism by countability and first-order sentences. (The prototype is  $\mathbb{Q}$ , characterized by Cantor's theorem, as we saw.) In the spirit of geometry (where, since Klein's Erlanger Programm, we have known of a connection between axiomatizability and symmetry), the following remarkable result was found independently by Engeler (1959), Ryll-Nardzewski (1959) and Svenonius (1959):

**Theorem 2.1.** *The countable structure  $M$  is  $\aleph_0$ -categorical if and only if  $\text{Aut}(M)$  is oligomorphic.*

The proof involves the well-known tool of “back-and-forth”. Without going into details (familiar to the experts), I note that back-and-forth is used, not just to show that two countable structures are isomorphic, but also to describe orbits on  $n$ -tuples structurally. The pre-requisite for back-and-forth is the ability to extend finite isomorphisms one point at a time.

A countable relational structure  $M$  is homogeneous if every isomorphism between finite substructures of  $M$  can be extended to an automorphism of  $M$ . The age of  $M$ , written  $\text{Age}(M)$ , is the class of finite structures embeddable in  $M$ . Now back-and-forth shows that  $M$  is homogeneous if and only if, for any  $X, Y \in \text{Age}(M)$  with  $X \subseteq Y$ , any embedding of  $X$  into  $M$  can be extended to  $Y$ . It suffices to require this when  $|Y| = |X| + 1$ . Using these ideas, Fraïssé (1953) showed:

**Theorem 2.2.** *A class  $\mathcal{A}$  of finite structures is the age of a countable homogeneous structure  $M$  if and only if  $\mathcal{A}$  is closed under isomorphism and under taking substructures, has only countably many non-isomorphic members, and has the amalgamation property. If these conditions hold, then  $M$  is unique up to isomorphism.*

We say that  $M$  is the *Fraïssé limit* of the class  $\mathcal{A}$ . This result is extremely useful for constructing examples. For instance, the class of finite triangle-free graphs has Fraïssé's properties; so there is a unique countable homogeneous universal triangle-free graph  $T$  (Henson 1971).

There is a natural topology on the symmetric group  $S_\infty$ , that of pointwise convergence. In this topology, the full automorphism group of any first-order structure is closed in  $S_\infty$ . Moreover, for  $G \leq H$ ,  $G$  is a dense subgroup of  $H$  if and only if  $G$  and  $H$  have the same orbits on  $n$ -tuples for all  $n$ . Now, given any permutation group  $G$ , there is a “canonical relational structure”  $M$  such that  $G$  is a dense subgroup of  $\text{Aut}(M)$ , and  $M$  is homogeneous.

The automorphism group of a countable homogenous structure over a finite relational language is thus a closed oligomorphic permutation group. The converse is false, however; I do not know any nice characterization of this class of permutation groups. (The closed oligomorphic groups are precisely the automorphism groups of countable  $\aleph_0$ -categorical structures; we may add the word “homogeneous” to the right-hand side of this equivalence.)

If  $M$  is a homogeneous structure and  $G = \text{Aut}(M)$  is oligomorphic, then  $f_n(G)$  is equal to the number of unlabelled  $n$ -element structures in the class  $\text{Age}(M)$  (i.e. up to isomorphism), and  $F_n(G)$  to the number of labelled structures (i.e. on the point set  $\{0, \dots, n-1\}$ ). Thus, enumeration of unlabelled or labelled structures in classes satisfying Fraïssé's hypotheses is equivalent to description of the appropriate orbit-counting sequence for an oligomorphic group  $G$ . As hinted above, many interesting combinatorial enumeration problems are of this type.

### 3. Generating Functions

It is common practice in combinatorial enumeration to use different forms of generating function in labelled and unlabelled enumeration problems (Goulden and Jackson 1983). Thus, we define the *ordinary generating function*

$$f_G(t) = \sum_{n \geq 0} f_n(G) t^n$$

for  $(f_n(G))$ , and the *exponential generating function*

$$F_G(t) = \sum_{n \geq 0} F_n(G) t^n / n!$$

for  $(F_n(G))$ . Convergence properties of these functions connect with growth rates (e.g. finite non-zero radius of convergence of  $f_G$  is equivalent to exponential growth of  $(f_n(G))$ ), but usually we treat the series formally. Both turn out to be specializations of a series in infinitely many variables, as follows.

If  $H$  is a finite permutation group of degree  $n$ , its *cycle index* is the polynomial

$$Z(H; s_1, \dots, s_n) = \frac{1}{|H|} \sum_{h \in H} s_1^{c_1(h)} \dots s_n^{c_n(h)},$$

where  $c_d(h)$  is the number of  $d$ -cycles in the cycle decomposition of  $h$ . Its rôle in enumeration is well known. Now, if  $G$  is a finite or oligomorphic permutation group, we define the *modified cycle index* of  $G$  to be

$$\tilde{Z}(G; s_1, s_2, \dots) = \sum_i Z(H_i; s_1, s_2, \dots),$$

where the summation is over representatives of the  $G$ -orbits on finite sets, and  $H_i$  is the group induced on the  $i$ -th set by its setwise stabilizer. (By convention, the empty set contributes a term 1.) For finite  $G$ , it can be shown that

$$\tilde{Z}(G; s_1, s_2, \dots) = Z(G; s_1 + 1, s_2 + 1, \dots).$$

A couple of examples for infinite groups are

$$\tilde{Z}(S_\infty) = \exp\left(-\sum_{n \geq 1} \frac{s_n}{n}\right)$$

and

$$\tilde{Z}(A) = \frac{1}{1 - s_1}.$$

These series specialize as follows:

**Proposition 3.1.** *For any oligomorphic permutation group  $G$ ,*

- (a)  $f_G(t) = \tilde{Z}(G; t, t^2, t^3, \dots)$ ;
- (b)  $F_G(t) = \tilde{Z}(G; t, 0, 0, \dots)$ .

The behaviour under direct and wreath products and point stabilizers can be described:

**Proposition 3.2.** (a)  $\tilde{Z}(G \times H) = \tilde{Z}(G)\tilde{Z}(H)$ ;

$$(b) \tilde{Z}(G \text{ Wr } H) = \tilde{Z}(H; \tilde{Z}(G) - 1);$$

$$(c) \sum_i \tilde{Z}(G_{\alpha_i}) = \frac{\partial}{\partial s_1} \tilde{Z}(G).$$

(The substitution in (b) is defined by

$$A(B) = A(B(s_1, s_2, s_3, \dots), B(s_2, s_4, s_6, \dots), B(s_3, s_6, s_9, \dots), \dots).$$

In (c), the summation is over a set of orbit representatives  $\alpha_i$ ;  $G_{\alpha_i}$  is the stabilizer of  $\alpha_i$ , acting on the remaining points.)

Thus,  $f_{G \text{ Wr } H}$  can be determined from  $f_G$  and  $\tilde{Z}(H)$ . For example,

$$f_{G \text{ Wr } S_\infty}(t) = \prod_{n \geq 1} (1 - t^n)^{-f_n(G)},$$

$$f_{G \text{ Wr } A}(t) = \frac{1}{2 - f_G(t)}.$$

I close this section with a few simple examples.

**Example 1.**  $f_{C_2}(t) = 1 + t + t^2$ ; so  $f_{C_2 \text{ Wr } A}(t) = 1/(1 - t - t^2)$ . This is the generating function for the Fibonacci numbers.

**Example 2.**  $f_n(S_\infty \text{ Wr } S_\infty)$  is the partition function  $p(n)$ , whose generating function is  $\prod_{n \geq 1} (1 - t^n)^{-1}$ . Moreover, the generating function for  $f_n(S_\infty \text{ Wr } S_\infty \text{ Wr } S_\infty)$  is  $\prod_{n \geq 1} (1 - t^n)^{-p(n)}$ , which converges for all  $t$ ; so the growth rate is slower than

exponential, though faster than  $\exp(n^{1-\varepsilon})$  for any  $\varepsilon > 0$ . This sequence arises in work of Cayley (1889) counting canonical forms.

**Example 3.** We have  $F_{G \text{ Wr } H}(t) = F_H(F_G(t) - 1)$ , where the usual substitution is intended. In particular,

$$F_{S_\infty \text{ Wr } G}(t) = F_G(e^t - 1).$$

Let  $F^*(t)$  be the exponential generating function for  $F_n^*(G)$ . Since  $F_n^*(G) = F_n(S_\infty \text{ Wr } G)$ , we have

$$F_G^*(t) = F_G(e^t - 1),$$

which is thus equivalent to the relation involving Stirling numbers given in Section 1.

**Example 4.** Let  $C$  be the group preserving the cyclic order on the roots of unity. Then

$$\tilde{Z}(C) = 1 + \sum_{n \geq 1} \frac{1}{n} \sum_{d|n} \phi(d) s_d^{n/d}.$$

From the fact that  $f_C(t) = 1/(1-t)$ , it can be deduced that

$$\exp(t/(1-t)) = \prod_{n \geq 1} (1 - t^n)^{-\phi(n)/n}.$$

For several further amusing examples, including groups realizing several familiar sequences, see Cameron (1987b), (1989). There are also close connections with Joyal's (1981) combinatorial formal power series. (Joyal regards an age as a category whose morphisms are the embeddings; the objects of cardinality  $n$  occur as the coefficients in formal power series.)

#### 4. The Random Graph

It follows from Fraïssé's theorem (2.2) that there is a unique countable homogeneous graph  $R$  in which every finite graph is embedded.  $R$  is characterized by the property that, if  $U$  and  $V$  are finite disjoint sets of vertices, there is a vertex  $z$  joined to every vertex in  $U$  and to no vertex in  $V$ . (This property translates the equivalent of homogeneity given just before the statement of (2.2).)

The graph  $R$  first appeared in the literature in a paper by Erdős and Rényi (1963), who showed that, with probability 1, a countable random graph is isomorphic to  $R$ . (The probability measure is defined by the rule "choose edges independently with probability  $\frac{1}{2}$ "; but in fact the same graph is obtained if we take any fixed edge probability  $p$  with  $0 < p < 1$ , or even if we let  $p$  vary a bit, e.g. tend to infinity not too slowly.) The proof of this paradoxical assertion is remarkably easy. First, given fixed finite disjoint sets  $U$  and  $V$ , the probability that the required vertex exists is 1. Now, since there are only countably many choices for  $U$  and  $V$ , and a countable intersection of sets of measure 1 has measure 1, the characteristic property of  $R$  holds with probability 1.

Although an explicit description of  $R$  is unnecessary for this discussion (the probabilistic argument shows its existence, while back-and-forth gives uniqueness), there are a couple of simple constructions for it:

*Construction 1.* The vertices are the natural numbers;  $x$  is joined to  $y$  if  $x < y$  and  $2^x$  occurs in the (binary) expression for  $y$  as a sum of distinct powers of 2 (or vice versa).

*Construction 2.* The vertices are the primes congruent to 1 (mod 4);  $p$  and  $q$  are joined if  $p$  is a quadratic residue mod  $q$ . (This is symmetric, by quadratic reciprocity.)

In Construction 1, the asymmetric form of the relation gives a model for the Zermelo-Fraenkel axioms of set theory excluding the axiom of infinity. The proof of the characteristic property of  $R$  in Construction 2 is a pleasant exercise using the Chinese Remainder Theorem and Dirichlet's Theorem.

As noted, the back-and-forth argument shows not only the uniqueness of  $R$ , but the homogeneous action of its automorphism group. If we use the first explicit construction of  $R$ , we find a group of primitive recursive automorphisms acting homogeneously on  $R$ . It would be interesting to investigate this group; for example, to see how the “recursive presentation” of  $R$  affects the structure of the group.

Truss (1985) showed that the full automorphism group of  $R$  is simple, and described all the cycle types of its elements. We know that  $F_n(\text{Aut}(R))$  and  $f_n(\text{Aut}(R))$  are equal to the numbers of labelled and unlabelled graphs on  $n$  vertices respectively; the former is  $2^{\frac{1}{2}n(n-1)}$ , the latter asymptotically  $2^{\frac{1}{2}n(n-1)}/n!$ .

$R$  is not the only countable homogeneous graph. All such graphs were found by Lachlan and Woodrow (1980) (the finite ones had been found earlier by Gardiner (1976)).

**Theorem 4.1.** *A countably infinite homogeneous graph is one of the following:*

- (i) *the disjoint union of  $m$  complete graphs of size  $n$ , where at least one of  $m$  and  $n$  is infinite;*
- (ii) *complement of (i) (complete multipartite);*
- (iii) *the Fraïssé limit of the class of finite graphs containing no complete subgraph of size  $n$  ( $n \geq 3$ );*
- (iv) *complement of (iii);*
- (v) *the random graph  $R$ .*

The universal  $K_n$ -free graphs in (iii) were constructed by Henson (1971). We met Henson's graph for  $n = 3$  earlier, where it was called  $T$ .

I turn now to more general relational structures. Suppose that we have some class  $\mathcal{X}$  of objects, each of which is described by a countable sequence of choices. There are two ways of assigning structure to the set  $\mathcal{X}$ :

*Method 1.* By assigning non-negative numbers summing to 1 to the outcomes of each choice,  $\mathcal{X}$  becomes a measure space. For example, a countable graph can be determined by choosing “edge” or “non-edge” for each pair of vertices; if we give each choice the value  $\frac{1}{2}$ , we obtain the above model. Equivalently, we could assign  $1/2^n$  to each possible extension of an  $n$ -vertex graph to an  $(n+1)$ -vertex graph. This technique can in principle be extended to many other classes of relational structures; but it is not clear what values to assign in general.

*Method 2.* We can make  $\mathcal{X}$  into a complete metric space by defining the distance between two objects to be a suitable (decreasing) functions of the length of the common initial sequence of choices defining them. (The longer we have to wait to distinguish two objects, the closer they are.) Although the metric involves a choice of function, the topology does not. A set  $\mathcal{Y} \subseteq \mathcal{X}$  is open if, for any  $Y \in \mathcal{Y}$ , there is an initial subsequence of the choice sequence defining  $Y$  which forces membership in  $\mathcal{Y}$ ; and a set  $\mathcal{Y} \subseteq \mathcal{X}$  is dense if, following any finite number of choices, there is a continuation defining a member of  $\mathcal{Y}$ .

A set is *residual* if it contains a countable intersection of open dense sets. The *Baire category theorem* asserts that, in a complete metric space, a residual set is non-empty (and, perforce, dense). Residual sets are regarded as “large”, comparable to sets of measure 1 in a probability space.

Analogously to the Erdős-Rényi theorem, it is possible to show that the set of graphs isomorphic to  $R$  is residual in the set of countable graphs. We now generalize this observation.

Let  $\mathcal{A}$  be any age. We define  $\mathcal{X}(\mathcal{A})$  to be the class of all structures on the set  $\mathbb{N}$  whose age is contained in  $\mathcal{A}$ . (Thus, for example, if  $\mathcal{A} = \text{Age}(R)$ , then  $\mathcal{X}(\mathcal{A})$  is the set of all graphs on the vertex set  $\mathbb{N}$ .) An element of  $\mathcal{A}$  is determined by countably many choices, the  $n^{\text{th}}$  choice describing how to extend an element of  $\mathcal{A}$  on the set  $\{0, \dots, n-1\}$  to one on the set  $\{0, \dots, n\}$ . So  $\mathcal{X}(\mathcal{A})$  is a complete metric space. Now we have:

**Proposition 4.2.** *Let  $M$  be a countable homogeneous structure. Then the set of structures isomorphic to  $M$  is residual in  $\mathcal{X}(\text{Age}(M))$ .*

No such result holds for measure, which seems much harder to deal with. However, there are some specific results. For example,  $\mathbb{Q}$  is the random total order (where the measure is defined by making all possible orderings of any finite set equally likely).

## 5. Growth Rates

Quite a bit is known about possible growth rates of the sequence  $(f_n(G))$ . (Of course, these general results apply to the numbers of unlabelled structures in a class satisfying Fraïssé's conditions. However, some of them are known to hold for any age.) A basic fact is that this sequence is non-decreasing: see Cameron (1976). Most of the results are due to Pouzet (1981) and Macpherson (1985a), (1985b).

Pouzet showed that the rate of growth is either polynomial ( $an^d \leq f_n(G) \leq bn^d$ , where  $n \in \mathbb{N}$  and  $a, b > 0$ ) or faster than any polynomial. In the latter case, Macpherson found a fractional exponential lower bound  $\exp(n^{\frac{1}{2}-\varepsilon})$ , comparable to the partition function. For primitive groups, Macpherson's result is much more striking:

**Theorem 5.1.** *There is an absolute constant  $c > 1$  such that, if  $G$  is primitive, then either  $f_n(G) = 1$  for all  $n$ , or  $f_n(G) \geq c^n$  for all sufficiently large  $n$ .*

Macpherson gave  $c = 2^{\frac{1}{3}} - \varepsilon$ ; it is conjectured that the result holds with  $c = 2 - \varepsilon$  (this would be best possible, see below).

Polynomial growth of degree  $k - 1$  is realized by, among others,  $S_\infty^k$  (acting with  $k$  orbits) and  $S_\infty \text{Wr} S_k$  (acting imprimitively).  $S_\infty \text{Wr} S_\infty$  realizes the partition function. Other fractional exponential growth rates, roughly  $\exp(n^{\frac{p+1}{p+2}})$  for  $p \in N$ , can also be realized. We saw that the growth rate for  $S_\infty \text{Wr} S_\infty \text{Wr} S_\infty$  is faster than fractional exponential but slower than exponential. There are many imprimitive examples with exponential growth; we saw the Fibonacci numbers realized by  $C_2 \text{Wr} A$ .

Primitive groups exhibiting exponential growth are fairly rare. Most of them are automorphism groups of “treelike objects” (Cameron 1987b) related to the  $\mathbb{Q}$ -trees of combinatorial group theory (Alperin and Bass 1987). There are also structures related to circular orders, including Lachlan’s (1984) “circular tournament”  $L$ . The group of order preserving and reversing permutations of  $L$  has the slowest known growth rate of any primitive group ( $f_n(\text{Aut}(L)) \sim 2^{n-2}/n$ ).

For growth rates faster than exponential, we see in nature a gap between factorial growth (for the homogeneous pair of linear orders, we have  $f_n(G) = n!$ ), and growth like  $\exp(cn^2)$  (realized by the random graph, and projective and affine spaces over finite fields). A result of Macpherson (1987) throws some light on this. The *independence property* (Shelah 1978) forces growth at least  $\exp(cn^2)$ ; for homogeneous structures over finite languages, negating the independence property bounds the growth by  $\exp(n^{1+\varepsilon})$ . Also, for  $\omega$ -stable structures, the same gap occurs, the criterion being the types of *strictly minimal* sets around which the structure is built. (See Cherlin–Harrington–Lachlan (1985) for the theory of  $\omega$ -stable,  $\aleph_0$ -categorical structures.)

In general, there is no upper bound for the rate of growth of  $(f_n(G))$ : take a homogeneous structure over a language where the number of  $n$ -ary relation symbols grows as fast as you please with  $n$ . On the other hand, for a homogeneous structure over a finite language,  $f_n(G) \leq \exp(P(n))$  for some polynomial  $P$ .

Many of the most interesting open questions about growth rate concern its smoothness. For example, do the limits

$$\begin{aligned} \lim_{n \rightarrow \infty} f_n(G)/n^d &\quad (\text{for polynomial growth of degree } d), \\ \lim_{n \rightarrow \infty} \log \log f_n(G) / \log n &\quad (\text{for fractional exponential growth}), \\ \lim_{n \rightarrow \infty} f_n(G)^{1/n} &\quad (\text{for exponential growth}), \end{aligned}$$

exist? If so, what possible values can these limits take? Apart from Macpherson’s gap in values of the third limit for primitive groups, almost nothing is known.

It is known that  $f_n(G) = f_{n+1}(G) = f_{n+2}(G)$  can only hold in the trivial case where  $G$  fixes a set of size at most  $n$  and is transitive on  $(n+2)$ -subsets of the complement. The situation  $f_n(G) = f_{n+1}(G)$  has been studied; a few examples are known, and some have been characterized, but a general result seems difficult.

## 6. Cyclic Automorphisms of Graphs

Measure and Baire category have been used for constructing subgroups of oligomorphic groups. Though much more general results are available, I will consider here only a special case, regular cyclic subgroups.

Let  $g$  be a cyclic automorphism of a countable graph  $\Gamma$ . Then the vertices of  $\Gamma$  can be indexed by the integers in such a way that  $g$  acts as a shift  $\alpha_i \mapsto \alpha_{i+1}$ . Let  $S = S(\Gamma, g) = \{n > 0 : \alpha_0 \sim \alpha_n\}$ . Then  $S$  determines  $\Gamma$  up to isomorphism ( $\alpha_i \sim \alpha_j$  if and only if  $|i - j| \in S$ ), and  $g$  up to conjugacy in  $\text{Aut}(\Gamma)$ . We write  $\Gamma = \Gamma_S$ ,  $g = g_S$ . Now we can ask: for which sets  $S$  is  $\Gamma_S$  isomorphic to some interesting graph?

A subset  $S$  of  $\mathbb{N}$  is determined in an obvious way by infinitely many choices; so the methods of measure and Baire category apply.

$S$  is called *universal* if every finite sequence of zeros and ones occurs as a consecutive subsequence of the characteristic function of  $S$ . It is easily checked that  $\Gamma_S \cong R$  if and only if  $S$  is universal. A weak form of the law of large numbers says that the set of universal sequences has measure 1; it can also be shown to be residual. Since residual sets and measure-1 sets have cardinality  $2^{\aleph_0}$ , we conclude that the random graph has  $2^{\aleph_0}$  non-conjugate cyclic automorphisms!

Now consider the homogeneous universal triangle-free graph  $T$ . First note that, for  $S \subseteq \mathbb{N}$ ,  $\Gamma_S$  is triangle-free if and only if  $S$  is *sum-free*, i.e.  $x, y \in S \Rightarrow x + y \notin S$ . Now sum-free sets can be determined by countably many binary choices, in an obvious way: considering natural numbers in turn, if  $n = x + y$  where  $x, y$  have already been put into  $S$ , then  $n \notin S$ ; otherwise we are free to choose.

Let  $S$  be a sum-free set. The only obvious necessary condition for a finite zero-one sequence  $\varepsilon$  to be a subsequence of the characteristic function of  $S$  is that, if  $j - i \in S$ , then  $\varepsilon_i$  and  $\varepsilon_j$  cannot both be 1. Call a sum-free set *sf-universal* if every finite zero-one sequence satisfying this condition is a subsequence of the characteristic function of  $S$ . Now  $\Gamma_S \cong T$  if and only if  $S$  is sf-universal.

Henson (1971) showed that  $T$  has cyclic automorphisms, i.e. that sf-universal sets exist. Can we prove this using category or measure? It is easily seen that a residual subset of all sum-free sets are sf-universal, so we do indeed get  $2^{\aleph_0}$  non-conjugate cyclic automorphisms of  $T$ . (Incidentally, I *conjecture* that an sf-universal set has density 0. This would, if true, give a “density version” of Schur’s theorem (1916) that  $\mathbb{N}$  cannot be covered by finitely many sum-free sets.)

However, when we turn to measure, the situation is different. Any set of odd numbers is obviously sum-free, and no such set is sf-universal. It came as a surprise to me to find that the probability that a random sum-free set consists entirely of odd numbers is non-zero. (This probability is about 0.218...). Furthermore, there are infinitely many “periodic” sum-free sets whose subsets have positive probability. (After the odd numbers, the next two are the congruence classes 2 and 3 (mod 5) and the congruence classes 1 and 4 (mod 5).) I do not know whether or not almost all sum-free sets are contained in periodic ones. (For details, see Cameron (1987a).)

Experiment suggests the possibility of quasi-periodic behaviour: a periodicity is established for a few cycles and then disrupted, to return with its phase shifted as part of a longer period. Maybe this process can continue infinitely and yield non-periodic sets whose subsets have positive probability.

What of the corresponding graphs  $\Gamma_S$ ? For each periodic set  $S$  whose subsets have positive probability, there is an “almost homogeneous” graph  $\Gamma^*$  such that  $\Gamma_{S'} \cong \Gamma^*$  for almost all subsets  $S'$  of  $S$ . (For the set of odd numbers,  $\Gamma^*$  is the “universal bipartite graph”.) It is not known whether any other triangle-free graphs occur with positive probability. (In particular, this is not known for  $T$ .)

Many other questions about random sum-free sets remain open. For example, what is the average density, and how is the density distributed? (There are “spectral lines” corresponding to the periodic sets described above, e.g. a delta-function of weight 0.218... at density  $\frac{1}{4}$  (from the sets of odd numbers). Does almost every sum-free set have a density? Is the density almost surely in the interval  $(0, \frac{1}{4}]$ ? Does the spectrum have a continuous part?)

I conclude with an example with very different behaviour. Covington (1989) calls a graph  $N$ -free if it contains no induced path of length 3. She showed that there is an “almost homogeneous” universal countable  $N$ -free graph  $C$ , unique up to isomorphism, but admitting no cyclic automorphisms. One can write down a condition on sets  $S$  equivalent to  $\Gamma_S$  being  $N$ -free. There are  $2^{\aleph_0}$  sets satisfying this condition, but the corresponding graphs  $\Gamma_S$  are pairwise non-isomorphic! In other words, an  $N$ -free graph has at most one conjugacy class of cyclic automorphisms.

## References

- Alperin, R., Bass, H. (1987): Length functions of group actions on  $A$ -trees. In: Combinatorial group theory and topology (Alta, Utah, 1984) (Ann. Math. Stud. 111). Princeton Univ. Press, Princeton, NJ, pp. 265–378
- Cameron, P.J. (1976): Transitivity of permutation groups on unordered sets. Math. Z. **148**, 127–139
- Cameron, P.J. (1987a): Portrait of a typical sum-free set. In: Whitehead, C. (ed), Surveys in combinatorics. London Math. Soc. (Lecture Notes, vol. 123). Cambridge Univ. Press, Cambridge, pp. 13–42
- Cameron, P.J. (1987b): Some treelike objects. Quart. J. Math. Oxford (2) **38**, 155–183
- Cameron, P.J. (1989): Some sequences of integers. Discrete Math. **75**, 85–102
- Cameron, P.J. (1990): Oligomorphic Permutation Groups. London Math. Soc. (Lecture Notes, vol. 152). Cambridge Univ. Press, Cambridge
- Cameron, P.J., Taylor, D.E. (1985): Stirling numbers and affine equivalence. Ars Combinatoria **20B**, 2–14
- Cantor, G. (1895): Beiträge zur Begründung der transfiniten Menge. Math. Ann. **46**, 481–512
- Cayley, A. (1889): Collected mathematical papers. Cambridge Univ. Press, London
- Cherlin, G.L., Harrington, L.A., Lachlan, A.H. (1985):  $\aleph_0$ -categorical,  $\aleph_0$ -stable structures. Ann. Pure Appl. Logic **28**, 103–135
- Covington, J. (1989): A universal structure for  $N$ -free graphs. Proc. London Math. Soc. (3) **58**, 1–16
- Engeler, E. (1959): Äquivalenzklassen von  $n$ -Tupeln. Z. Math. Logik Grundl. Math. **5**, 340–345
- Erdős, P., Rényi, A. (1963): Asymmetric graphs. Acta Math. Acad. Sci. Hungar. **14**, 295–315
- Fraïssé, R. (1953): Sur certains relations qui généralisent l’ordre des nombres rationnels. C. R. Acad. Sci. Paris **237**, 540–542
- Gardiner, A.D. (1976): Homogeneous graphs. J. Comb. Theory (B) **20**, 94–102

- Goulden, I.P., Jackson, D.M. (1983): Combinatorial enumeration. Wiley, New York
- Henson, C.W. (1971): A family of countable homogeneous graphs. *Pacific J. Math.* **38**, 69–83
- Joyal, A. (1981): Une théorie combinatoire des séries formelles. *Adv. Math.* **42**, 1–82
- Lachlan, A.H. (1984): Countable homogeneous tournaments. *Trans. Amer. Math. Soc.* **284**, 431–461
- Lachlan, A.H., Woodrow, R.E. (1980): Countable ultrahomogeneous undirected graphs. *Trans. Amer. Math. Soc.* **262**, 51–94
- Macpherson, H.D. (1985a): Orbits of infinite permutation groups. *Proc. London Math. Soc.* (3) **51**, 246–284
- Macpherson, H.D. (1985b): Growth rates in infinite graphs and permutation groups. *Proc. London Math. Soc.* (3) **51**, 285–294
- Macpherson, H.D. (1987): Permutation groups of rapid growth. *J. London Math. Soc.* (2) **35**, 276–286
- Pouzet, M. (1981): Application de la notion de relation presque-enchaînable au dénombrément des restrictions finies d'une relation. *Z. Math. Logik Grundl. Math.* **27**, 289–332
- Ryll-Nardzewski, C. (1959): On category in power  $\leq \aleph_0$ . *Bull. Acad. Pol. Sér. Math. Astr. Phys.* **7**, 545–548
- Schur, I. (1916): Über die Kongruenz  $x^m + y^m \equiv z^m \pmod{p}$ . *Jber. Deutsch. Math.-Verein.* **25**, 114–117
- Shelah, S. (1978): Classification theory and the number of non-isomorphic models. North-Holland, Amsterdam
- Svenonius, L. (1959):  $\aleph_0$ -categoricity in first-order predicate calculus. *Theoria* **25**, 82–94
- Truss, J.K. (1985): The group of the countable universal graph. *Math. Proc. Camb. Philos. Soc.* **98**, 213–245



# Geometric Presentations of Groups with an Application to the Monster

Alexander A. Ivanov

Institute for System Studies, Academy of Sciences of the USSR, 9, Prospect 60 Let Oktyabrya, 117312, Moscow, USSR

## 1. Introduction

Let  $\mathcal{G} = (\mathcal{G}, I, \Delta, t)$  be a *geometry*, that is a set  $\mathcal{G}$  of elements together with a symmetric reflexive *incidence relation*  $I$  and a *type function*  $t : \mathcal{G} \rightarrow \Delta$ . A geometry is supposed to satisfy the following axiom: *the restriction of  $t$  to any maximal set of pairwise incident elements is a bijection onto  $\Delta$ .*

Let  $\Phi$  be a *flag* of  $\mathcal{G}$ , i.e. a set of pairwise incident elements. Let  $\mathcal{G}_\Phi = (\mathcal{G}_\Phi, I_\Phi, \Delta_\Phi, t_\Phi)$  where  $\mathcal{G}_\Phi$  is the set of elements which are not contained in  $\Phi$  and are incident to all elements in  $\Phi$ ;  $\Delta_\Phi = \Delta - t(\Phi)$ ;  $I_\Phi$  and  $t_\Phi$  are the restrictions on  $\mathcal{G}_\Phi$  of  $I$  and  $t$ , respectively. Then  $\mathcal{G}_\Phi$  is a geometry in the above sense and is called the *residual geometry* of  $\mathcal{G}$  with respect to  $\Phi$ .

Let  $\mathcal{G}$  and  $\mathcal{G}'$  be geometries over the same set of types. A mapping  $\phi : \mathcal{G}' \rightarrow \mathcal{G}$  is a *morphism of geometries* if  $\phi$  preserves the incidence relation and the type function. A morphism is called a *covering* if its restriction to any proper residual geometry is an isomorphism. A morphism from a geometry onto itself is an *automorphism*. Let  $G \leq \text{Aut}(\mathcal{G})$  be an automorphism group of  $\mathcal{G}$ .  $G$  is said to be *flag-transitive* if it acts transitively on the set of maximal flags of  $\mathcal{G}$ .

We usually assume that  $\Delta = \{0, 1, \dots, r - 1\}$  where  $r$  is the *rank* of the geometry. A geometry  $\mathcal{G}$  is *connected* if the graph with  $\mathcal{G}$  as vertices and  $I$  as edges is connected. All geometries we will consider are supposed to be connected. A considerable amount of information about geometry is carried by its *diagram*. The latter is a graph on  $\Delta$  where the edge joining  $i$  and  $j$  symbolizes the rank 2 residual geometries of type  $\{i, j\}$ . The empty edge stands for the generalized digons (any two elements of different types are incident), an ordinary edge stands for projective planes, etc.

Let  $\mathcal{G}$  be a geometry,  $G$  be a flag-transitive group of automorphisms of  $\mathcal{G}$  and  $\Phi = \{\alpha_0, \alpha_1, \dots, \alpha_{r-1}\}$  be a maximal flag of  $\mathcal{G}$ . Let  $G_i = G(\alpha_i)$  be the stabilizer of  $\alpha_i$  in  $G$  and  $\mathcal{A}$  be the amalgam of these stabilizers. The members  $G_i$  of  $\mathcal{A}$  are the *maximal parabolics*. In the flag-transitive case  $\mathcal{G}$  can be reconstructed from  $G$  and  $\mathcal{A}$ , namely its elements of type  $i$  are the (right) cosets of  $G_i$  in  $G$ ,  $0 \leq i \leq r - 1$ ; two cosets are incident if they have a nonempty intersection. In this situation we write  $\mathcal{G} \cong \mathcal{G}(G, \mathcal{A})$ .

A geometry  $\mathcal{G}$  possesses a *universal covering*  $\phi_u : \tilde{\mathcal{G}} \rightarrow \mathcal{G}$  such that for any other covering  $\phi : \mathcal{G}' \rightarrow \mathcal{G}$  there is a covering  $\psi : \tilde{\mathcal{G}} \rightarrow \mathcal{G}'$  such that  $\phi_u = \phi\psi$ . Let  $G$  act flag-transitively on  $\mathcal{G}$ . Then automorphisms from  $G$  can be lifted to automorphisms of

$\tilde{\mathcal{G}}$  and all these liftings form a group  $\tilde{G}$  which acts flag-transitively on  $\tilde{\mathcal{G}}$ . Let  $\mathcal{A}$  be an amalgam. By definition an  $\mathcal{A}$ -group is a group which contains  $\mathcal{A}$  and is generated by the elements of  $\mathcal{A}$ . An  $\mathcal{A}$ -homomorphism is a homomorphism of  $\mathcal{A}$ -groups whose restriction on  $\mathcal{A}$  is the identity mapping. If the class of  $\mathcal{A}$ -groups is nonempty then there exists a universal  $\mathcal{A}$ -group  $U(\mathcal{A})$  such that any  $\mathcal{A}$ -group is an image of  $U(\mathcal{A})$  under an  $\mathcal{A}$ -homomorphism. One can define  $U(\mathcal{A})$  as a group having all elements of  $\mathcal{A}$  as generators and all equalities valid in the members of  $\mathcal{A}$  as relations.

The following fundamental result was proved almost simultaneously in [Pas1, Tit2] and in an unpublished manuscript by S.V. Shpectorov.

**Theorem 1.1.** *Let  $\mathcal{G}$  be a geometry,  $G$  be a flag-transitive automorphism group of  $\mathcal{G}$  and  $\mathcal{A}$  be the amalgam of maximal parabolics. Let  $\phi_u : \tilde{\mathcal{G}} \rightarrow \mathcal{G}$  be the universal covering,  $\tilde{G}$  be the group of all liftings of the automorphisms from  $G$  and  $U(\mathcal{A})$  be the universal  $\mathcal{A}$ -group. Then  $\tilde{\mathcal{G}} \cong \mathcal{G}(U(\mathcal{A}), \mathcal{A})$  and  $\tilde{G} \cong U(\mathcal{A})$ .*

A geometry whose universal covering is an isomorphism is said to be *simply connected*. By Theorem 1.1  $\mathcal{G}$  is simply connected if and only if  $G$  is the universal  $\mathcal{A}$ -group.

In terms of generators and relations Theorem 1.1 gives a practical method for the determination of universal coverings of flag-transitive geometries. In this method one writes down a presentation for  $U(\mathcal{A})$  starting with  $\mathcal{A}$  and constructs  $U(\mathcal{A})$  by a coset enumeration. Nowadays this method is rather popular and on its base a number of very interesting results on classification of flag-transitive geometries have been obtained (cf. [Hei, KT, Pas2, WY] and many other papers). In some papers generators with relations are constructed just from conditions on the geometry. The language of generators and relations is rather convenient for concrete calculations. On the other hand the language of amalgams is more abstract. The correlation between these languages looks like the correlation between the languages of matrices and modules in the theory of linear operators.

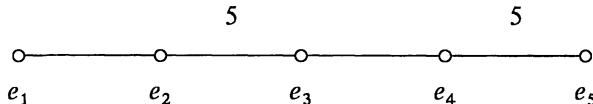
Let  $\mathcal{G}, G$  and  $\mathcal{A}$  be as above. A *geometric presentation* associated with the action of  $G$  on  $\mathcal{G}$  is a set of elements of  $\mathcal{A}$  and a set of relations on these elements valid in the members of  $\mathcal{A}$ , giving a presentation of  $U(\mathcal{A})$ . Notice that a particular geometric presentation is involved in the definition of  $U(\mathcal{A})$ . If the geometry  $\mathcal{G}$  is simply connected then  $G \cong U(\mathcal{A})$  and the geometric presentation is said to be *faithful*, so it is a presentation of  $G$ .

Construction of faithful geometric presentations for large sporadic groups form the content of the proposed lecture. In the next section we start with a simple illustrative example. We consider a presentation for Janko's group  $J_1$  from [ATLAS]. This presentation can be treated as a geometric one and it is faithful since the corresponding geometry is simply connected. In Sect. 3 we discuss practical methods for proving simple connectedness. Here we present two new ideas. One of them relies on consideration of simply connected subgeometries, the second one deals with the fixed points of involutions from the automorphism group. Section 4 is devoted to simple connectedness of geometries of the sporadic simple groups  $J_4$  (the Janko's group),  $F_2$  (the Baby Monster) and  $F_1$  (the Monster). The result for  $J_4$  implies

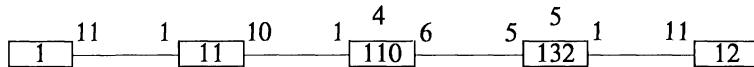
that a presentation from [SW] is faithful. In the case of  $F_1$  our result enabled S. Norton to prove that the famous Y-presentation (c.f. [ATLAS], [CNS]) is faithful. In Sect. 5 we discuss some uniqueness aspects of sporadic groups.

## 2. An Illustrative Example: The Group $J_1$

In [ATLAS] Janko's group  $J_1$  is presented as a Coxeter group having diagram



with the additional relations:  $z_1 = z_2 = z_3 e_1 = 1$ . Here  $z_i$  is the nontrivial element in the center of  $\langle e_i, e_{i+1}, e_{i+2} \rangle \cong A_5 \times Z_2$ ,  $1 \leq i \leq 3$ . Let  $H_i = \langle e_j | 1 \leq j \leq 5, j \neq i \rangle$ ,  $2 \leq i \leq 5$  be subgroups of  $H \cong J_1$ . Then the relations for  $H$  which involve only generators of  $H_i$  give a presentation for  $H_i$  and  $H_2 \cong H_4 \cong A_5 \times Z_2$ ,  $H_3 \cong D_6 \times D_{10}$ ,  $H_5 \cong L_2(11)$ . Let  $\mathcal{H}$  be the amalgam of the subgroups  $H_i$ ,  $2 \leq i \leq 5$  and  $\mathcal{G}(J_1) = \mathcal{G}(H, \mathcal{H})$ . The geometry  $\mathcal{G}(J_1)$  possesses a nice combinatorial interpretation. The group  $J_1$  acts distance-transitively on a graph  $\Gamma(J_1)$  of valency 11 with the following intersection diagram:



Then the elements of type 2, 3, 4 and 5 in  $\mathcal{G}(J_1)$  are the Petersen subgraphs, the length 5 cycles fixed by  $Z_3$ -subgroups, edges and vertices, respectively. The incidence corresponds to the symmetrized inclusion.

The above presentation is a geometric one with respect to the action of  $J_1$  on  $\mathcal{G}(J_1)$ . So it is faithful if and only if  $\mathcal{G}(J_1)$  is simply connected. Let  $\phi: \tilde{\mathcal{G}} \rightarrow \mathcal{G}(J_1)$  be the universal covering. Then  $\phi$  induces a covering  $\hat{\phi}: \tilde{\Gamma} \rightarrow \Gamma(J_1)$  of graphs. Each 5-cycle in  $\Gamma(J_1)$  is contained in a Petersen subgraph. Thus to prove that  $\hat{\phi}$  is an isomorphism one should split all cycles of  $\Gamma(J_1)$  into cycles of length 5. The existence of these splittings is proved in [Ivn2, Lemma 6.9] by induction on the length of cycles.

The above presentation was independently obtained in [Ivn1] just as a consequence of the simple connectedness of  $\mathcal{G}(J_1)$ .

## 3. Proving Simple Connectedness of Geometries

A direct combinatorial proof of simple connectedness of a geometry basically consists of the following two steps (cf. [IS2, Pas2, Ron2]).

**Step 1.** With the geometry  $\mathcal{G}$  in question one associates a graph  $\Gamma = \Gamma(\mathcal{G})$  with the property that the universal covering  $\phi: \tilde{\mathcal{G}} \rightarrow \mathcal{G}$  induces a covering  $\hat{\phi}: \tilde{\Gamma} \rightarrow \Gamma$  of

graphs. From the condition that  $\phi$  is a covering of  $\mathcal{G}$  one makes a conclusion that cycles of  $\Gamma$  from a certain class  $\mathcal{K}$  are *contractible* with respect to  $\hat{\phi}$ . This means that each cycle from  $\mathcal{K}$  can be lifted to an isomorphic cycle in  $\tilde{\Gamma}$ .

**Step 2.** Now to prove that  $\hat{\phi}$  (and hence  $\phi$ ) is an isomorphism it is sufficient to show that the normal closure of the cycles from  $\mathcal{K}$  generates the fundamental group of  $\Gamma$ . This is equivalent to the claim that an arbitrary cycle in  $\Gamma$  can be split into cycles from  $\mathcal{K}$ .

There is a number of ways to associate a graph with a geometry. The typical strategy is the following. One chooses two types in  $\mathcal{G}$  (points and lines) and considers the *point graph* of  $\mathcal{G}$ , that is, the graph on the set of points where adjacency corresponds to collinearity. Here we consider a special situation when  $\mathcal{G}$  contains a simply connected subgeometry.

Let  $\mathcal{G}, G$  and  $\mathcal{A}$  be as above and let  $\mathcal{F}$  be a simply connected subgeometry in  $\mathcal{G}$  over the set of types  $0 \in \mathcal{A}_{\mathcal{F}} \subseteq \mathcal{A}$ . Suppose that the stabilizer  $F$  of  $\mathcal{F}$  in  $G$  acts flag-transitively on  $\mathcal{F}$ . Let  $\mathcal{B}$  be the subamalgam of  $\mathcal{A}$  consisting of the subgroups  $F_i = G_i \cap F$ ,  $i \in \mathcal{A}_{\mathcal{F}}$ . Then by Theorem 1.1 the simple connectedness of  $\mathcal{F}$  implies that  $\mathcal{B}$  generates in  $\tilde{G} \cong U(\mathcal{A})$  a subgroup isomorphic to  $F$  and it is easy to see that this subgroup stabilizes in  $\tilde{\mathcal{G}}$  a subgeometry which is an isomorphic lifting of  $\mathcal{F}$ . Let  $\mathcal{C}$  be the set of all images of  $\mathcal{F}$  under  $G$ . We assume that distinct subgeometries from  $\mathcal{C}$  have distinct sets of elements of type 0. Let  $\mathcal{C}(\alpha_0)$  be the set of subgeometries from  $\mathcal{C}$  which contain  $\alpha_0$ . Then  $G_0$  acts on  $\mathcal{C}(\alpha_0)$  as it acts on the cosets of  $F_0$ . Let  $R(\alpha_0)$  be a nontrivial symmetric 2-orbit in this action (if any such exists). Let  $\Sigma$  be a graph on  $\mathcal{C}$  where two subgeometries are adjacent if they have an element  $\beta$  of type 0 in common and are in the relation  $R(\beta)$ .  $\Sigma$  will be called an *intersection graph of subgeometries*. Let  $\mathcal{F}$  and  $\mathcal{F}'$  be adjacent in  $\Sigma$ . Let  $m = m(\mathcal{G})$  be the number of elements  $\beta$  of type 0 in  $\mathcal{F} \cap \mathcal{F}'$  such that  $\mathcal{F}$  and  $\mathcal{F}'$  are in the relation  $R(\beta)$ . Then the valency of  $\Sigma$  is equal to  $|\mathcal{F}^0| \cdot v/m$ , where  $\mathcal{F}^0$  is the set of elements of type 0 in  $\mathcal{F}$ , and we have the following

**Lemma 3.1.** *If  $m(\tilde{\mathcal{G}}) = m(\mathcal{G})$  then  $\phi$  induces a covering  $\hat{\phi}: \tilde{\Sigma} \rightarrow \Sigma$  where  $\tilde{\Sigma}$  is an intersection graph of subgeometries in  $\tilde{\mathcal{G}}$ .*

**Corollary 3.2.** *Let  $v = 1$  and suppose that  $F_0$  is maximal in  $F$ . Then the valency of  $\Sigma$  is equal to  $|\mathcal{F}^0|$  and  $\phi$  induces a covering  $\hat{\phi}: \tilde{\Sigma} \rightarrow \Sigma$ .*

Let  $\mathcal{E}(\alpha_0)$  be a subgraph in  $\Sigma$  having  $\mathcal{C}(\alpha_0)$  as vertices and  $R(\alpha_0)$  as edges. This subgraph can be lifted to an isomorphic subgraph in  $\tilde{\Sigma}$ . This implies that a cycle from  $\mathcal{E}(\alpha_0)$  as well as its images under  $G$  are contractible. An additional set of contractible cycles can give other classes of simply connected subgeometries in  $\mathcal{G}$  (if any).

The intersection graphs are crucial in our simple connectedness proof for the geometries in Sect. 4. The point is that for  $G \cong J_4$ ,  $F_2$  and  $F_1$  the corresponding geometry contains a simply connected subgeometry whose stabilizer is a maximal order subgroup of  $G$ . It turns out that in each case an intersection graph of

subgeometries corresponds to the minimal antireflexive 2-orbit of  $G$  acting on the subgeometries.

Now let us turn to the second step. In all examples we meet in the lecture,  $\mathcal{K}$  contains all shortest cycles of  $\Gamma$ . These are pentagons in the  $J_1$ -geometry and triangles in the geometries of  $J_4$ ,  $F_2$  and  $F_1$ . So we come to a *standard problem*: to present the cycles of  $\Gamma$  as sums of the shortest ones. If the shortest cycles are triangles then the standard problem is the problem of *triangulation* of graphs.

Let us give a sufficient condition for a graph to be triangulable (compare [Ron 2, Lemma 5]). As usual  $\Gamma_j(x)$  denotes the set of vertices of  $\Gamma$  which are at distance  $j$  from a vertex  $x$ ,  $0 \leq j \leq d$ , where  $d$  is the *diameter* of  $\Gamma$ .

**Lemma 3.3.** *If for each  $2 \leq j \leq d$  the conditions (i) and (ii) are satisfied then  $\Gamma$  is triangulable*

- (i) *For  $y \in \Gamma_j(x)$  the subgraph induced by  $\Gamma_1(x) \cap \Gamma_{j-1}(y)$  is connected.*
- (ii) *For  $y, z \in \Gamma_j(x)$ ,  $y$  and  $z$  adjacent, either  $\Gamma_1(x) \cap \Gamma_{j-1}(y) \cap \Gamma_{j-1}(z) \neq \emptyset$  or  $\Gamma_1(x) \cap \Gamma_{j-1}(y)$  and  $\Gamma_1(x) \cap \Gamma_{j-1}(z)$  are joined by an edge.*

Let  $G$  act on  $\Gamma$  and let  $\tau$  be an involution in  $G$ . Let  $\Theta(\tau)$  be the subgraph of  $\Gamma$  induced by the vertices fixed by  $\tau$  and  $\mathcal{K}(\tau)$  be the set of cycles from  $\mathcal{K}$  which lie in  $\Theta(\tau)$ . Then we can try to prove that normal closure of  $\mathcal{K}(\tau)$  generates the fundamental group of  $\Theta(\tau)$ . If we would succeed doing this for all representatives of conjugacy classes of involutions in  $G$  then we will get the following nice criterion: if a cycle is fixed by an involution then it is contractible.

## 4. Geometric Presentations of Large Sporadic Groups

### 4.1 Generic Properties

Let  $G$  be one of the groups  $J_4$ ,  $F_2$  or  $F_1$ . Then  $G$  contains an elementary abelian subgroup  $K$  of order  $2^r$  ( $r = 4$  for  $J_4$  and  $r = 5$  otherwise) such that  $N_G(K)/C_G(K) \cong L_r(2)$ . Let  $1 < K_0 < K_1 < \dots < K_{r-1} = K$  be a chain of subgroups in  $K$ . Let  $\mathcal{G}(G)$  be a geometry whose elements of type  $i$  are the subgroups of  $G$  conjugate to  $K_i$ ,  $0 \leq i \leq r-1$  and the incidence be defined by inclusion. Let  $G_i = N_G(K_i)$ ,  $0 \leq i \leq r-1$  and  $\mathcal{A}$  be the amalgam of these subgroups. Then  $\mathcal{A}$  is the amalgam of maximal parabolics and  $\mathcal{G}(G) = \mathcal{G}(G, \mathcal{A})$ . The diagram of  $\mathcal{G}(G)$  is a *string*. Namely the edge  $\{i, j\}$  is empty for  $|i - j| > 1$ . The residue of type  $\{i, i+1\}$ ,  $0 \leq i \leq r-3$  is the projective plane over  $GF(2)$ . For  $\mathcal{G}(J_4)$  and  $\mathcal{G}(F_2)$  the residue of type  $\{r-2, r-1\}$  is the geometry of edges and vertices of the Petersen graph while for  $\mathcal{G}(F_1)$  it is the famous triple cover of the generalized quadrangle of order  $(2, 2)$  denoted by  $\circ \overline{\phantom{x}} \circ$ .

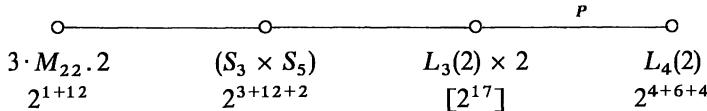
The parabolic  $G_0$  is of the form  $G_0 \cong A \cdot B$  where  $A = O_2(G_0)$  is an extraspecial group  $2_+^{1+n}$ ,  $n = 12, 22$  and  $24$ ;  $B \cong 3 \cdot M_{22} \cdot 2$ ,  $Co_2$  and  $Co_1$  for  $G \cong J_4$ ,  $F_2$  and  $F_1$ , respectively. The subgroup  $A$  coincides with the kernel of  $G_0$  acting on the residual geometry  $\mathcal{G}(\alpha_0)$  of  $\mathcal{G}$  with respect to  $\alpha_0$ . The geometry  $\mathcal{G}(\alpha_0)$  possesses a *natural representation* in the module  $\bar{A} = A/Z(A)$ . This means that elements of type  $i$  in

$\mathcal{G}(\alpha_0)$  correspond to certain  $i$ -dimensional subspaces of  $\bar{A}$ ,  $1 \leq i \leq r - 1$  and the incidence in  $\mathcal{G}(\alpha_0)$  corresponds to inclusion of the subspaces. In particular the image of  $K$  in  $\bar{A}$  corresponds to an element of type  $r - 1$ . This property can be taken as a definition of  $K$ .

The geometries  $\mathcal{G}(J_4)$  and  $\mathcal{G}(F_2)$  are so-called *P-geometries*. A characterization of a nice class of girth 5 graphs is reduced to the classification of *P-geometries* [Inv3]. The simple connectedness of  $\mathcal{G}(J_4)$  and  $\mathcal{G}(F_2)$  was proved just within this classification. Other results in this direction can be found in [IS1, IS2, Shp1, Shp2, Shp3]. The geometry  $\mathcal{G}(F_1)$  and the truncations of  $\mathcal{G}(J_4)$  and  $\mathcal{G}(F_2)$  over type  $r - 1$  are the minimal 2-local parabolic geometries of these groups described in [RSt].

#### 4.2 Janko's Group $J_4$

The action of  $J_4$  on  $\mathcal{G}(J_4)$  is described by the following *diagram of parabolics* where under the node of type  $i$  the structure of  $G_i$  is given ( $[2^n]$  stands for an arbitrary group of order  $2^n$ ).



The parabolic  $G_2 = N_G(K_2)$  is contained in a subgroup  $L \cong 2^{11} : M_{24}$ , and  $K_2 \leq O_2(L)$ . Let us consider the set of subgroups of  $O_2(L)$  conjugated in  $L$  to nontrivial subgroups of  $K_2$ . Then we obtain a subgeometry  $\mathcal{F}$  in  $\mathcal{G}(J_4)$  over the type set  $\{0, 1, 2\}$ . This subgeometry is isomorphic to the minimal parabolic geometry  $\mathcal{G}(M_{24})$  of  $M_{24}$ . Notice that  $\mathcal{G}(M_{24})$  arises as a residue of type  $\{2, 3, 4\}$  in  $\mathcal{G}(F_1)$ .

The following result was proved in [Hei] and independently checked on a computer by S.V. Shpectorov and the author.

**Lemma 4.1.** *The geometry  $\mathcal{G}(M_{24})$  is simply connected.*

Thus we obtain a family  $\mathcal{C}$  of simply connected subgeometries in  $\mathcal{G}(J_4)$ .  $G_0$  acts on  $\mathcal{C}$  as  $M_{22} \cdot 2$  acts on the set of 77 blocks of  $S(3, 6, 22)$ . So there are two symmetric antireflexive 2-orbits in this action, having valencies 16 and 60. Let  $R(\alpha_0)$  be the 2-orbit of valency 60.

**Lemma 4.2.**  $m(\tilde{\mathcal{G}}) = m(\mathcal{G}(J_4)) = 7$ .

Hence by Lemmas 3.1 and 4.2 the universal covering of  $\mathcal{G}(J_4)$  induces a covering of an intersection graph of subgeometries  $\Sigma(J_4)$  of valency  $15, 180 = 1771 \cdot 60/7$ . Notice that the elements of type 0 in  $\mathcal{G}(M_{24})$  are the 1771 sextets of  $S(5, 8, 24)$ .

The subgraph  $\Xi(\alpha_0)$  contains representatives of all classes of triangles. So for simple connectedness of  $\mathcal{G}(J_4)$  it is sufficient to prove that  $\Sigma(J_4)$  is triangulable. The latter fact was proved by application of Lemma 3.3 and an analysis of the subgraphs induced by the fixed points of involutions. A description of the 2-point stabilizers of  $J_4$  acting on the cosets of  $L \cong 2^{11} : M_{24}$  given in [Nor1] was used in a crucial way. A detailed proof is presented in [Ivn4].

The parabolic  $G_3$  is contained in a subgroup  $H \cong 2^{10} : L_5(2)$ . It is shown in [Ivn2] that the permutable action of  $G$  on the cosets of  $H$  preserves a graph  $\Pi$  of valency 31 and girth 5. Let  $\mathcal{H}(J_4)$  be a rank 3 geometry whose elements are vertices, edges and pentagons of  $\Pi$  with the natural incidence relation. In [SW] an explicit geometric presentation associated with the action of  $J_4$  on  $\mathcal{H}(J_4)$  is given. This is a compact and nice presentation which is similar to the Steinberg presentation for the group  $D_5(2)$ . Notice that  $H$  is isomorphic to a maximal parabolic in  $D_5(2)$ . As a corollary of the simple connectedness of  $\mathcal{G}(J_4)$  it was proved in [Ivn4] that  $\mathcal{H}(J_4)$  is simply connected.

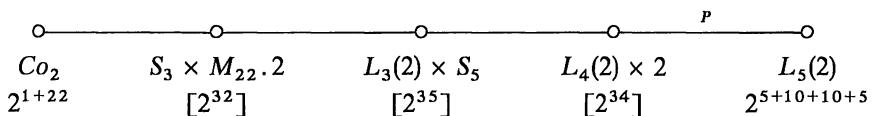
**Theorem 4.3.** *The geometries  $\mathcal{G}(J_4)$  and  $\mathcal{H}(J_4)$  are simply connected and the presentation by G. Stroth and R. Weiss is faithful.*

The above result implies also the simple connectedness of a rank 3 2-local geometry of  $J_4$  from [Bue] where maximal parabolics are  $L$ ,  $G_0$  and a subgroup  $2^{3+12}.(S_5 \times L_3(2))$ .

It is mentioned in [AS] that an independent proof of the triangulability of  $\Sigma(J_4)$  is obtained by the authors of that paper.

### 4.3 The Baby Monster $F_2$

The action of  $F_2$  on  $\mathcal{G}(F_2)$  is described by the following diagram of parabolics.



Let  $S \cong 2^{9+16}.Sp_8(2)$  and  $E \cong 2 \cdot 2E_6(2).2$  be subgroups of  $G \cong F_2$  and let  $\bar{S} = S/O_2(S)$ ,  $\bar{E} = E/O_2(E)$ . Then up to a suitable choice of  $S$  with respect to  $K \cong 2^5$ , the intersections  $S \cap G_i$ ,  $0 \leq i \leq 3$  contain  $O_2(S)$  and their images in  $\bar{S}$  are the maximal parabolics of the natural geometry of  $\bar{S}$ . Up to a suitable choice of  $E$  with respect to  $K$  and  $S$  the intersections  $E \cap G_i$ ,  $0 \leq i \leq 2$ ,  $E \cap S$  contain  $O_2(E)$  and their images in  $\bar{E}$  are the maximal parabolics of the natural geometry of  $\bar{E}$ .

So  $S$  gives us a subgeometry  $\mathcal{S}$  of type  $C_4$  while  $E$  gives a subgeometry  $\mathcal{E}$  which is a truncation of a geometry of type  $F_4$ . Both  $\mathcal{S}$  and  $\mathcal{E}$  are simply connected by [Tit1], [Tit2]. Let  $\mathcal{C}$  be the family of subgeometries which contains  $\mathcal{E}$ . Put  $E_0 = E \cap G_0$ . Then  $E_0$  is maximal in  $E$  and  $[O_2(G_0) : O_2(G_0) \cap E_0] = 2$ . So we can apply Corollary 3.2 and come to a covering of an intersection graph  $\Sigma(F_2)$  of subgeometries. Then  $\Sigma(F_2)$  is a graph on the class of  $\{3, 4\}$ -transpositions in  $F_2$  where two transpositions are adjacent if their product is a central involution of  $F_2$ . The subgroup  $S$  acting on  $\Sigma(F_2)$  has an orbit of length 120 which induces a complete subgraph  $\Theta$ . Since  $\mathcal{S}$  is simply connected one can show that  $\Theta$  is contractible. On the other hand  $\Theta$  contains representatives of all triangles in  $\Sigma(F_2)$ . So we come to the triangulation problem. Here we are in a position to prove the triangulability just by Lemma 3.3 using the structure constants of  $\Sigma(F_2)$  calculated in [Hig]. Eventually we come to the following result whose proof will be given in [Ivn5].

**Theorem 4.4.** *The geometry  $\mathcal{G}(F_2)$  is simply connected.*

It would be interesting to find a nice explicit presentation of  $F_2$  as a corollary of Theorem 4.4.

Some independent results on triangulability of  $\Sigma(F_2)$  were obtained by T. Meixner [Mei] in his study of a  $c$ -extension of the natural geometry of  ${}^2E_6(2)$  related to  $F_2$  (cf. [Bue]).

#### 4.4 The Monster $F_1$

In the case of  $\mathcal{G}(F_1)$  we have the following diagram

$$\begin{array}{ccccccc} \circ & \circ & \circ & \circ & \sim & \circ \\ Co_1 & S_3 \times M_{24} & L_3(2) \times 3.S_6 & L_4(2) \times S_3 & & L_5(2) \\ 2^{1+24} & [2^{35}] & [2^{39}] & [2^{39}] & & 2^{5+1+5+10+10+5} \end{array}$$

$\mathcal{G}(F_1)$  contains  $\mathcal{G}(F_2)$  as a subgeometry. This embedding can be described as follows.  $G \cong F_1$  contains exactly two conjugacy classes ( $2A$  and  $2B$ ) of involutions with centralizers  $2 \cdot F_2$  and  $2_+^{1+24} \cdot Co_1$ , respectively. Let  $K$  be the subgroup involved in the definition of  $\mathcal{G}(F_1)$  and  $\tau \in K^\#$ . Then  $\tau$  is of type  $2B$  and  $O_2(C_G(\tau))$  contains an involution  $\sigma$  of type  $2A$  such that  $C_G(\sigma) \cap N_G(K)$  induces  $L_5(2)$  on  $K$ . Let  $\mathcal{F}$  consist of all images of the subgroups of  $K$  under conjugation by  $C_G(\sigma)$ . Then  $\mathcal{F}$  is isomorphic to  $\mathcal{G}(F_2)$  and its stabilizer  $F$  is  $C_G(\sigma) \cong 2 \cdot F_2$ . So by Theorem 4.4 we obtain a family  $\mathcal{C}$  of simply connected subgeometries in  $\mathcal{G}(F_1)$  whose members are in a one-to-one correspondence with the  $2A$ -involutions in  $G$ .

The subgroup  $F_0 = F \cap G_0$  is maximal in  $F$  and  $[O_2(G_0) : O_2(G_0) \cap F_0] = 2$ . So by Corollary 3.2 the universal covering of  $\mathcal{G}(F_1)$  induces a covering  $\hat{\phi}$  of a graph  $\Sigma = \Sigma(F_1)$ . The latter is a graph on the set of  $2A$ -involutions in  $G$  where two involutions are adjacent if their product is a  $2B$ -involution. So to prove that  $\mathcal{G}(F_1)$  is simply connected it is sufficient to prove that the normal closure of the cycles from  $\mathcal{E}(\alpha_0)$  and their conjugacies under  $G$  generates the fundamental group of  $\Sigma$ . Let  $\Sigma^* = \Sigma^*(F_1)$  be the *Monster graph*, that is a graph on the same set as  $\Sigma$  where two involutions are adjacent if their product is an involution of type  $2A$ . Detailed information about  $\Sigma^*$  is contained in [Nor2], [GMS]. Let us show that  $\hat{\phi}$  induces a covering of  $\Sigma^*$ . Let  $x, y \in 2A$  and  $x \cdot y \in 2A$ . Then  $y \in \Sigma_2(x)$ . Let  $\Omega$  be a graph on the set  $\Sigma_1(x) \cap \Sigma_1(y)$  where  $u$  and  $v$  are adjacent if the triangle  $\{x, u, v\}$  is conjugated to a triangle from  $\mathcal{E}(\alpha_0)$ . Then  $\Omega$  is connected and hence  $\hat{\phi}$  induces a covering of  $\Sigma^*$ . Now,  $\Sigma^*$  is also an intersection graph of subgeometries. This means that two subgeometries from  $\mathcal{C}$  which are adjacent in  $\Sigma^*$  have a common element of type 0. Let  $\mathcal{E}^*(\alpha_0)$  be the subgraph in  $\Sigma^*$  on the set of subgeometries passing through  $\alpha_0$ . Then  $\mathcal{E}^*(\alpha_0)$  contains representatives of all triangles in  $\Sigma^*$ . The triangulability of  $\Sigma^*$  was proved by application of Lemma 3.3 with an analysis of the fixed points of involutions and using some results from [Nor2] and [GMS]. So we come to the following result whose proof will appear in [Ivn6].

**Theorem 4.5.** *The geometry  $\mathcal{G}(F_1)$  is simply connected.*

An independent triangulation proof for the Monster graph is given in [AS, Sect. 8].

Let  $\mathcal{A}$  be the amalgam of maximal parabolics associated with the action of  $F_1$  on  $\mathcal{G}(F_1)$ . Then by Theorems 1.1 and 4.5,  $U(\mathcal{A}) \cong F_1$ . Let  $\mathcal{B}$  be the subamalgam of  $\mathcal{A}$  consisting of the parabolics  $G_0$ ,  $G_1$  and  $G_2$ . The residues of types  $\{0, 1, 2\}$  and  $\{0, 1, 2, 3\}$  are classical and simply connected. Hence, if  $i = 4$  or  $5$  then  $G_i$  coincides with the universal group of the amalgam consisting of the subgroups  $G_i \cap G_j$  for  $0 \leq j \leq i - 1$ . So we have the following

**Corollary 4.6.**  $U(\mathcal{B}) \cong F_1$ .

By arguments similar to the above ones it is possible to prove simple connectedness of a geometry  $\mathcal{G}(Co_1)$  isomorphic to the residue  $\mathcal{G}(\alpha_0)$  in the Monster geometry. Here the crucial subgeometry is  $\mathcal{G}(Co_2)$  isomorphic to the analogous residue in the Baby Monster geometry. It is proved in [Shp3] that  $\mathcal{G}(Co_2)$  is simply connected. So in view of Lemma 4.1 and Theorem 4.5 the following proposition holds.

**Theorem 4.7.** *The geometry  $\mathcal{G}(F_1)$  is 2-simply connected.*

A geometry is 2-simply connected if any of its covering whose restrictions on rank 2 residues are isomorphisms is an isomorphism itself. To describe the universal 2-cover of the Baby Monster geometry  $\mathcal{G}(F_2)$  is an interesting open problem. I conjecture that the automorphism group of this 2-cover is a nonsplit extension of the elementary abelian group of order  $3^{4371}$  by  $F_2$ .

The above results were announced at the Durham Symposium on Groups and Combinatorics in July 1990. During this symposium S. Norton [Nor4], using his recent results [Nor3] and Corollary 4.6, proved that the famous Y-presentation (cf. [ATLAS, CNS]) is faithful. Namely he proved that the group  $Y_{555}$  in the notation of [ATLAS] is isomorphic to the Bimonster, i.e. to the wreath product of the Monster and a group of order 2. This result gives an explicit presentation for the Monster and for other related groups.

## 5. On Uniqueness of Sporadic Groups

Here we indicate some relationship between the above results and the uniqueness problem for sporadic simple groups.

We consider a geometric approach to the uniqueness problem which involves the following steps. (1) Starting with some general properties of a group  $G$  in question, prove that it should act flag-transitively on a geometry  $\mathcal{G}(G)$ . (2) Using geometrical properties of  $\mathcal{G}(G)$ , prove uniqueness of the corresponding amalgam of maximal parabolics. (3) Prove that  $\mathcal{G}(G)$  is simply connected.

Notice that this scheme is very close to the approach based on so-called *uniqueness systems* proposed in [AS].

First of all the existence of the geometries  $\mathcal{G}(J_4)$ ,  $\mathcal{H}(J_4)$ ,  $\mathcal{G}(F_2)$  and  $\mathcal{G}(F_1)$  follows from some general properties of the groups such as the structure of the centralizers

of involutions. In the case of  $J_4$  the second step was realized in [SW] for  $\mathcal{H}(J_4)$  in terms of generators and relations and in [Shp] for  $\mathcal{G}(J_4)$  in terms of amalgams. So Theorem 4.3 implies the following

**Corollary 4.8.** *There is a unique group of type  $J_4$ .*

For  $F_2$  and  $F_1$  the step (2) is not done yet but we believe that [GMS] and [Sev] contain enough information for this purpose.

**Notes added in proof.** (1) The simple connectedness of the geometry  $\mathcal{G}(\text{Co}_1)$  which is a residue in the Monster geometry is proved in [Ivanov, A.A.: The minimal parabolic geometry of the Conway group  $\text{Co}_1$  is simply connected. In: Proc. Int. Conf. "Combinatorics 90", Gaeta, Italy 1990. (To appear)].

(2) The geometric approach to the uniqueness problem (cf. Section 5) is now completed for the groups  $F_1$  and  $F_2$  [Ivanov, A.A.: A geometric approach to the uniqueness problem for the sporadic simple groups. Dokl. Akad. Nauk SSSR **316** (1991) 1043–1046 (Russian)].

(3) Using the simple connectedness of the Baby Monster  $P$ -geometry the author has proved that  $Y_{433}$  is isomorphic to  $2 \times 2 \cdot F_2$ . This gives an explicit presentation for  $F_2$  and hence answers the problem posed in the paragraph after Theorem 4.4 [Ivanov A.A.: Presenting the Baby Monster. Preprint, 1991].

## References

- [AS] Aschbacher, M., Segev, Y.: Extending morphisms of groups and graphs. (To appear)
- [Bue] Buekenhout, F.: Diagram geometries for sporadic groups. Contemp. Math. **45** (1985) 1–31
- [ATLAS] Conway, J.H., et al.: Atlas of Finite Groups. Oxford Univ. Press, 1985
- [CNS] Conway, J.H., Norton, S.P., Soicher, L.H.: The Bimonster, the group  $Y_{555}$  and the projective plane of order 3. In: Tangora, M.C.(ed.) Computers in algebra. Marcel Dekker, 1988, pp. 27–50
- [GMS] Griess, R.J., Jr., Meierfrankenfeld, U., Segev, Y.: A uniqueness proof for the Monster. Ann. Math. **130** (1989) 567–602
- [Hei] Heiss, S.: On a parabolic system of type  $M_{24}$ . J. Algebra (to appear)
- [Hig] Higman, D.G.: A monomial character of Fischer's Baby Monster. In: Gross, F. (ed.) Proc. of the Conference on Finite Groups. Acad. Press, New York, 1976, pp. 277–283
- [Ivn1] Ivanov, A.A.: A geometric characterization of the group  $J_1$ . In: Proc. X All-Union Conf. on Group Theory. Gomel', 1986, p. 91 (Russian)
- [Ivn2] Ivanov, A.A.: On 2-transitive graphs of girth 5. Europ. J. Combin. **8** (1987) 393–420
- [Ivn3] Ivanov, A.A.: Graphs of girth 5 and diagram geometries related to the Petersen graph. Dokl. Akad. Nauk. **295** (1987) 529–533 (Russian) [English transl. Sov. Math. Dokl. **36** (1988) 83–87]
- [Ivn4] Ivanov, A.A.: A presentation for  $J_4$ . Proc. London Math. Soc. (To appear)
- [Ivn5] Ivanov, A.A.: A geometric characterization of Fischer's Baby Monster. University of Western Australia, Preprint, May 1991.

- [Ivn6] Ivanov, A.A.: A geometric characterization of the Monster. Proc. Symp. "Groups and Combinatorics", Durham 1990. (To appear)
- [IS1] Ivanov, A.A., Shpectorov, S.V.: Geometries for sporadic groups related to the Petersen graph. II. Europ. J. Combin. **10** (1989) 347–361
- [IS2] Ivanov, A.A., Shpectorov, S.V.: The  $P$ -geometry for  $M_{23}$  has no nontrivial 2-coverings. Europ. J. Combin. **11** (1990) 373–379
- [KTs] Komissartschik, E.A., Tsaranov, S.V.: Construction of finite group amalgams and geometries. Geometries of the group  $U_4(2)$ . Comm. Algebra **18** (1990) 1071–1118
- [Mei] Meixner, T.: Personal communication, 1990
- [Nor1] Norton, S.: The construction of  $J_4$ . In: Cooperstein, B. and Mason, G. (eds.) The Santa Cruz Conference on Finite Groups. Amer. Math. Soc., 1980, pp. 271–278
- [Nor2] Norton, S.: The uniqueness of the Fischer-Griess Monster. Contemp. Math. **45** (1985) 271–285
- [Nor3] Norton, S.: Presenting the monster? Bull. Math. Soc. of Belgium (A) **42** (1990) 595–605
- [Nor4] Norton, S.: Constructing the Monster. Proc. Symp. "Groups and Combinatorics", Durham 1990. (To appear)
- [Pas1] Pasini, A.: Some remarks on covers and apartments. In: Baker, C.A. and Batten, L.M. (eds.) Finite geometries. Marcel Dekker Inc., New York Basel 1985, pp. 233–250
- [Pas2] Pasini, A.: A classification of a class of Buekenhout geometries exploiting amalgams and simple connectedness. Preprint 1989
- [Ron1] Ronan, M.: Covers and automorphisms of chamber systems. Europ. J. Combin. **1** (1980) 259–269
- [Ron2] Ronan, M.: Coverings of certain finite geometries. In: Cameron, P., Hirschfeld, J., and Hughes, D. (eds.) Finite geometries and designs. Cambridge Univ. Press, 1981, pp. 316–331
- [RSt] Ronan, M.A., Stroth, G.: Minimal parabolic geometries for the sporadic groups. Europ. J. Combin. **5** (1984) 59–91
- [Seg] Segev, Y.: On the uniqueness of Fischer's baby monster (to appear)
- [Shp1] Shpectorov, S.V.: A geometric characterization of the group  $M_{22}$ . In: Klin, M.H., Faradzev, I.A. (eds.) Investigations in algebraic theory of combinatorial objects. Moscow, VNIISI, 1985, pp. 112–123 (Russian)
- [Shp2] Shpectorov, S.V.: On geometries with the diagram  $P''$ . Preprint, 1988 (Russian)
- [Shp3] Shpectorov, S.V.: The universal 2-cover of the  $P$ -geometry  $\mathcal{G}(Co_2)$  (To appear)
- [SW] Stroth, G., Weiss, R.: Modified Steinberg relations for the group  $J_4$ . Geom. Dedic. **25** (1988) 513–525
- [Tit1] Tits, J.: A local approach to buildings. In: The Geometric Vein. The Coxeter Festschrift. Springer, Berlin 1981, pp. 519–547
- [Tit2] Tits, J.: Ensembles ordonnés, immeubles et sommes amalgamées. Bull. Soc. Math. Belg. **A38** (1986) 367–387
- [WY] Weiss, R., Yoshiara, S.: A geometric characterization of the groups Suz and HS. J. Algebra **133** (1990) 182–196



# Some Developments in Ramsey Theory

Vojtech Rödl

Department of Mathematics, Emory University, Atlanta, GA 30322, USA

## Introduction

Ramsey Theory is a part of combinatorial mathematics that studies the behaviour of structures under partitions. Many Ramsey type results assert that complete disorder is impossible – they find some regular substructures in general combinatorial structures.

Generic results are due to F.P. Ramsey [R3] and B.L. van der Waerden [V]. The systematic study of these and related statements was initiated by the work of P. Erdős and R. Rado. This includes extending van der Waerden's theorem by giving a complete description of those systems of linear equations that are preserved under partitions [R1, R2], developing the partition calculus – a theory which deals with transfinite extensions of Ramsey's theorem [EHR, EHMR] – and setting many further directions for systematic study [Er1, E4]. Another step was done by Hales and Jewett [HJ], who found a new Ramsey type combinatorial principle, the theorem for partitions of words over finite alphabets. This result was further generalized by R.L. Graham and B.L. Rothschild [GR2] for partitions of multi-dimensional words (parameter words). This subsequently led R.L. Graham, K. Leeb and B.L. Rothschild [GLR] to the proof of the conjecture of G. C. Rota on finite vector spaces. Moreover, the Hales-Jewett Theorem played an important role in studying induced subgraphs, distributive lattices and  $(m, p, c)$ -sets. A survey of the work in this area can be found in the book “Ramsey Theory” by R.L. Graham, B.L. Rothschild and J. Spencer [GRS].

Most of the work in Ramsey theory has been focused toward two directions: structural Ramsey type theorems (finding new statements of Ramsey type) and quantitative Ramsey type theorems (evaluating the size and number of the corresponding objects). Results were found that deal with vector spaces (and their combinatorial analogues – parameter words), set and relational-systems, groups and lattices, and other structures. Various generalizations were found for partitions restricted by topological concepts. (cf. [CS])

Ramsey type statements have a number of striking applications and connections to other parts of mathematics (e.g. [PH, CW, AH, KV] cf. also [NR11]). On the other hand some deep theorems were discovered using non-combinatorial means (e.g. [Fu, FK1, and FK2]).

In this account, we will outline the development in a few areas of this subject.

## The Ramsey Theorem and Ramsey Numbers

We will adopt the usual convention of identifying the positive integer  $m$  with the set of its predecessors  $\{0, 1, 2, \dots, m - 1\}$ . For a set  $X$ ,  $[X]^k$  denotes the set of all  $k$ -element subsets of  $X$ .

**The Ramsey Theorem (1930).** *For all positive integers  $n_1, n_2, \dots, n_l$  and  $k$  there exists an integer  $m$  such that for every partition*

$$[m]^k = C_1 \cup C_2 \cup \dots \cup C_l$$

*there exists  $i$  and a set  $Y_i \subseteq m = \{0, 1, 2, \dots, m - 1\} \mid |Y_i| = n_i$  such that  $[Y_i]^k \subseteq C_i$ .*

The Ramsey number  $R_k(n_1, n_2, \dots, n_l)$  is the smallest integer  $m$  such that the Ramsey Theorem is valid. Clearly  $R_1(n_1, n_2, \dots, n_l) = \sum_{i=1}^l (n_i - 1) + 1$ ; to determine the Ramsey numbers for  $k \geq 2$  is much harder. Even in the next simplest case, i.e.,  $k = 2, l = 2$  progress on determining the corresponding Ramsey numbers  $R(n_1, n_2) = R_2(n_1, n_2)$  has been very slow. For example, the lower bound for  $R(n, n)$  due to P. Erdős [E1] was already known in 1947:

$$R(n, n) > (1 + o(1)) \frac{1}{e\sqrt{2}} n 2^{n/2}. \quad (1)$$

In more than forty years since its proof, this bound has been improved only by a factor of 2. This was done by J. Spencer [S4] implementing the Lovász Local Lemma [EL]. The upper bound

$$R(n_1, n_2) \leq \binom{n_1 + n_2 - 2}{n_1 - 1} \quad (2)$$

given by Erdős and Szekeres [ES] in 1935 has been improved for  $n_1$  fixed and  $n_2 \rightarrow \infty$  by Graver and Yackel [GY] and Ajtai, Komlós and Szemerédi [AKS]. For other values of  $n_1, n_2$  there have been improvements on (1) only by a constant factor ([Fr1]). The bound in (2) was improved in [R6]. (see also [GR1])

$$R(n_1, n_2) \leq C_1 \left( \frac{n_1 + n_2}{n_1} \right) / [\log(n_1 + n_2)]^{c_2}. \quad (3)$$

Subsequently, for values of  $n_1$  and  $n_2$  which are close to each other ( $\frac{n_2}{\sqrt{\log n_2}} \ll n_1 \leq n_2$ ), (3) was improved by A. Thomason [T2] who established

$$R(n_1, n_2) \leq C_3 \left( \frac{n_1 + n_2}{n_1} \right) \exp \left[ -\frac{n_1}{2n_2} \log n_2 + C\sqrt{\log n_2} \right]. \quad (4)$$

The proof of both inequalities is based on the idea of counting triangles in graphs and their complements (cf. [Go, L2]). Both (3) and (4) can probably be significantly improved.

Of all asymptotic bounds for  $R(n_1, n_2)$  the only satisfactory ones are known for  $n_1 = 3$  and  $n_2 = n$  large:

$$C_1 \left( \frac{n}{\log n} \right)^2 \leq R(3, n) \leq C_2 \frac{n^2}{\log n}.$$

The lower bound was established by P. Erdős [E2] and the upper bound by Ajtai, Komlós, Szemerédi [AKS] (cf. Shearer [S2]).

For  $n_1 \geq 4$  and  $n_2 = n$  large, the situation is much less satisfactory. For instance,

$$C_3 \left( \frac{n}{\log n} \right)^{5/2} \leq R(4, n) \leq C_4 \frac{n^3}{\log^2 n}$$

(cf. [S5] [AKS]) are the best bounds currently known. A substantial effort has gone into finding the exact values  $R(n_1, n_2)$  for small values of  $n_1, n_2$ . It is unlikely that  $R(n_1, n_2)$  will be determined for many new values of  $n_1$  and  $n_2$  (cf. table below). For example, it would appear that determination of  $R(5, 5)$  will require some essentially new ideas. The reader is referred to [GRS], [RK] and [R10] for further discussion.

$n_1/n_2$	3	4	5	6	7	8	9
3	6	9	14	18	23	28	36
4		18	25–27	34–43	47–66		
5			43–52	51–94	76–160		
6				102–169			

Let us also mention that various extensions of the Ramsey numbers have been investigated quite extensively in the last 15 years. For a survey on this work, see [GR1]. We will close this section by mentioning three old problems.

1) From the bound on Ramsey numbers  $R(n, n)$  one can deduce that  $\sqrt{2} \leq R(n, n)^{1/n} \leq 4$ .

So, the first problem is to improve either of these bounds for  $n \rightarrow \infty$ . (i.e. improve the base of the exponent in either case) It is, for example, not even known whether the  $\lim_{n \rightarrow \infty} R(n, n)^{1/n}$  exists.

2) Another old problem is to decide the growth of  $R_2(3, 3, \dots, 3)$  ( $n$ -times). It is not even known whether  $\lim_{n \rightarrow \infty} [R_2(3, 3, \dots, 3)]^{1/n}$  is finite or infinite. (The existence of the limit here is easy to prove). The known bounds are

$$(315)^{1/5} \leq [R_2(3, 3, \dots, 3)]^{1/n} \leq \left( \frac{n}{e} \right) (1 + o(1)). \quad (5)$$

Note that the upper bound in (5) has not been improved in the 74 years since its discovery [S1].

3) Perhaps the most important problem concerning Ramsey numbers is to decide about the behavior of  $R_3(n, n)$ . The following is known (cf. [EHMR]) and it is believed that the upper bound is correct

$$\frac{n^2}{6}(1 + o(1)) \leq \log_2 R_3(n, n) \leq 2^{cn}.$$

## Structural Induced and Restricted Ramsey-Type Theorems

Most of my contributions to Ramsey Theory have been in this area, much of this is joint work with J. Nešetřil. The area was opened by work of J. Folkman [Fo]. Answering a question of Erdős and Hajnal, he proved that for every  $n$  there exists a graph  $H$  not containing  $K_{n+1}$ , but for which every 2-coloring of the edges results in a monochromatic  $K_n$ . A related result was obtained by W. Deuber [D1], P. Erdős, A. Hajnal, L. Posa [EHP] and myself [R4]. These results, together with work of Graham, Leeb, Rothschild and others, (cf. [D2, DR, L1, NR1]) layed the ground work for further systematic investigation of Ramsey classes.

Let  $\mathbf{K}$  be a class of objects (e.g. graphs, parameter sets, ...) endowed with isomorphism and subobjects. For  $A, B \in \mathbf{K}$ , consider the set  $\binom{B}{A}$  consisting of all subobjects of  $B$  which are isomorphic to  $A$  ( $A$ -subobjects of  $B$ ). Let  $A \in \mathbf{K}$ . The class  $\mathbf{K}$  has the  $A$ -Ramsey Property if for every  $B \in \mathbf{K}$  there exists  $C \in \mathbf{K}$  with  $C \rightarrow (B)_2^A$ , i.e. for every 2-coloring of the  $A$ -subobjects of  $C$  there exists a  $B$ -subobject of  $C$  with all its  $A$ -subobjects colored the same. If  $\mathbf{K}$  has the  $A$ -Ramsey property for each  $A \in \mathbf{K}$  we say that  $\mathbf{K}$  is a *Ramsey class*.

In this language, the classical Ramsey theorem means that the class of all finite sets together with inclusion forms a Ramsey class. Also the class of all finite dimensional vector spaces over finite fields is a Ramsey class – this is the statement of Graham-Leeb-Rothschild Theorem [GLR].

Effort has been focused toward finding new Ramsey classes and, for the other classes, one has tried to characterize  $A \in \mathbf{K}$  for which  $\mathbf{K}$  has the  $A$ -Ramsey property.

So, for example, the class of all  $k$ -uniform hypergraphs with linearly ordered vertex sets,  $\text{Gra}(k)$  (more general, the class  $\text{Soc}(\mathcal{A})$  of all set systems of type  $\mathcal{A}$ ) is a Ramsey class. This was proved in [NR3] (see also [NR5]) and independently by Abramson and Harrington [AH]. On the other hand, the same classes with vertex sets not ordered fail to be Ramsey. Here the  $A$ -Ramsey property holds only for highly symmetric hypergraphs  $A$ , i.e., those which are either complete or empty. Subsequently, analogous statements were proved for the class of partially ordered sets [NR6] and the classes of parameter sets and vector space systems, [P1, P2].

Another direction was to consider the same type of questions for certain sub classes  $\mathbf{K}$  of the class of all  $k$ -uniform hypergraphs, partially ordered sets, and others. We will explain here the result of [NR3].

A type  $\mathcal{A} = (m_\delta; \delta \in \mathcal{A})$  is an indexed collection of positive integers. A system  $A$  of type  $\mathcal{A}$  is a pair  $(X, \mathbf{M})$  where  $X$  is a finite linearly ordered set,  $\mathbf{M} = (\mathbf{M}_\delta; \delta \in \mathcal{A})$  and  $\mathbf{M}_\delta \in [X]^{m_\delta}$ . The elements of  $\mathbf{M}$  are called *edges*. The system  $A$  is a *subsystem*

(subobject) of  $B = (Y, \mathbf{N})$  if  $X$  is a subset of  $Y$  with the induced order, and  $\mathbf{M}_\delta = \mathbf{N}_\delta \cap P(X)$  for every  $\delta \in \Delta$ . Two systems  $(X, \mathbf{M})$  and  $(Y, \mathbf{N})$  are isomorphic if there is a monotone bijection  $f: X \rightarrow Y$  taking  $\mathbf{M}_\delta$  onto  $\mathbf{N}_\delta$  for every  $\delta \in \Delta$ . Let  $\binom{B}{A}$  be the set of all subsystems of  $B$  which are isomorphic to  $A$ . Let  $\text{Soc}(\Delta)$  be the class of all set systems of type  $\Delta$ . For example, if  $\Delta = (2)$  the class  $\text{Soc}(\Delta)$  is the class of graphs with ordered vertices and induced embedding.

A system  $F \in \text{Soc}(\Delta)$  is called irreducible if every pair of points of  $F$  is contained in an edge of  $F$ . For  $\mathbf{F} \subseteq \text{Soc}(\Delta)$  ( $\mathbf{F}$  may be infinite), denote by  $\text{Forb}(\mathbf{F})$  the set of all  $A \in \text{Soc}(\Delta)$  with  $\binom{A}{F} = \emptyset$  for every  $F \in \mathbf{F}$ .

The following theorem was proved in [NR3], (see also [NR5]).

Let  $\mathbf{F} \subseteq \text{Soc}(\Delta)$  be a (possibly infinite) set of irreducible set systems. Then  $\mathbf{K} = \text{Forb}(\mathbf{F})$  is a Ramsey class. Moreover, if  $\mathbf{K}$  satisfies certain additional assumptions, the implication can be reversed, i.e., if  $\mathbf{K} \subseteq \text{Soc}(\Delta)$  is a Ramsey class satisfying additional assumptions then  $\mathbf{K}$  is of the form  $\mathbf{K} = \text{Forb}(\mathbf{F})$ , where  $\mathbf{F}$  is the class of irreducible set systems (cf. also [N1]).

As there are many naturally defined classes that are not described in terms of forbidden irreducible set systems (note for example that for type  $\Delta = (2)$  the only irreducible systems are complete graphs) the following general problem (asked in many particular forms by P. Erdős) arises:

Given  $A \in \text{Soc}(\Delta)$  determine which classes  $\mathbf{K} \subseteq \text{Soc}(\Delta)$  have the  $A$ -Ramsey property.

The answer to this question is hard even in the simplest cases, i.e. for  $|A| = 1$  and  $\Delta = (2)$ . However, quite general satisfactory answers can be found.

In [NR4] and [NR6], partite amalgamation was introduced. This method was successful in [NR7] [NR8] and moreover allowed simplification of (difficult) proofs of some earlier statements. It was furthermore observed that it is possible to apply this proof technique to Hales Jewett cubes and to arithmetic progression. This motivated further research with P. Frankl and R.L. Graham, [FRG] where extending the earlier results of H.J. Prömel [P1, P2], we succeeded in applying partite amalgamation to vector spaces and parameter words. This paper was followed by work of J. Nesetril, H.J. Prömel, B. Voigt and myself ([NR9, NR10, PV]) which describes the general form of results which can be obtained by partite amalgamation.

## Euclidean Ramsey Theory

Hadwiger and Nelson [H1] posed the problem of determining the chromatic number  $\chi(n)$  of Euclidean space  $\mathbf{R}^n$  – the minimum integer  $r + 1$  such that for every decomposition  $\mathbf{R}^n = X_1 \cup X_2 \cup \dots \cup X_r$  and every positive real  $\alpha$  there are two points in  $X_i$  for some  $i$  with distance precisely  $\alpha$ . It was proven in [H1], [MM], [W] that  $4 \leq \chi(2) \leq 7$  and it follows moreover from [Fa] that under the assumption that the sets  $X_i$  are measurable, the lower bound can be replaced by 5. The best current bounds for  $\chi(n)$  are

$$(1.2)^n \leq \chi(n) \leq (3 + o(1))^n.$$

The lower bound is due to Frankl and Wilson [FW], the upper bound to Larman and Rogers [LR].

Generalizing the problem of Hadwiger and Nelson, the following concept was investigated in [EG]. A finite subset  $A \subset \mathbf{R}^d$  is called *Ramsey* if for every  $r$  there exists  $n = n(r, A)$  such that for every partition  $\mathbf{R}^n = X_1 \cup X_2 \cup \dots \cup X_r$ , there is some  $i$  and  $A' \subset X_i$  with  $A'$  congruent to  $A$ .

In a series of papers, Erdős et al. [EG] have investigated this property. They have shown that all Ramsey sets  $A$  are *spherical*, that is,  $A$  has a circumcenter, a point equidistant to all points of  $A$ . On the other hand, it is proved in [EG] that the vertex sets of rectangular parallelepipeds (and therefore all their subsets) are Ramsey.

The simplest sets that are spherical but can not be embedded into rectangular parallelepipeds are obtuse triangles.

In [FR3], (cf. also [FR2]) it is shown that all vertex sets of nondegenerate simplices are Ramsey. Moreover, for both simplices and rectangular parallelepipeds one can in fact choose  $n(r, A) = C(A) \log r$ , where  $C(A)$  is an appropriate positive constant. This can be further strengthened as follows:

Let  $S^m(\varrho) = \{(x_0, x_1, \dots, x_m) \in \mathbf{R}^{m+1} : x_0^2 + x_1^2 + \dots + x_m^2 = \varrho^2\}$  be the  $m$ -dimensional sphere of radius  $\varrho$ . In [Gr1], a configuration  $A$  is called *sphere-Ramsey* if for every  $r$  there exists  $m = m(A, r)$  and  $\varrho = \varrho(A)$  such that for every partition  $S_\varrho^m = X_1 \cup \dots \cup X_r$ , there exists  $A' \subset X_i$  with  $A'$  congruent to  $A$ . R.L. Graham [Gr1] proved that rectangular parallelepipeds are sphere-Ramsey and asked whether for a parallelepiped  $A$  with circumradius  $\alpha$  one can choose  $\varrho(A) = \alpha + \delta$ , where  $\delta > 0$  is arbitrarily small. A positive answer to this question for 2-point configurations, follows from a result of Lovász [L3]. In [FR3], a positive answer to Graham's equation is given in the following stronger form:

Call a set  $A$  *hyper-Ramsey* if for all  $\delta > 0$  there exist positive constants  $c = c(A, \delta)$ ,  $\varepsilon = \varepsilon(\varepsilon, \delta)$  and subsets  $X = X(m) \subset S^m(\varrho(A) + \delta)$  for  $m > m_0(\delta)$  such that

- i)  $|X(m)| < c^m$  but if
- ii)  $|Y| \geq (1 - \varepsilon)^m |X(m)|$  and  $Y \subset X(m)$

then  $Y$  contains a congruent copy of  $A$ .

It has been stated in [FR3] that if  $A \subset \mathbf{R}^n$ ,  $B \subset \mathbf{R}^m$  are hyper-Ramsey, then, so is their product  $A * B = \{x * y | x \in A, y \in B\}$  where for  $x = (x_1, x_2, \dots, x_n)$  and  $y = (y_1, \dots, y_m)$  one defines  $x * y = (x_1, \dots, x_n, y_1, \dots, y_m)$ . From [FW], it follows that every 2-element set is hyper-Ramsey (cf. also [RS]). Thus, also the products of 2-element sets, i.e. all rectangular parallelepipeds are hyper-Ramsey (which answers Graham's conjecture affirmatively).

Our current knowledge does not exclude the possibility that all spherical sets are hyper-Ramsey, but we can not prove this even for non-degenerate simplices. The first step towards settling this problem would be to decide whether all obtuse triangles are hyper-Ramsey. We can, however, prove that non-degenerate simplices (and also their products) are exponentially sphere-Ramsey. This means that they are sphere-Ramsey and there exists  $\varepsilon = \varepsilon(A) > 0$  and constant  $C = C(A)$  such that  $m = m(A, r)$  can be chosen to satisfy  $m \leq C(1 + \varepsilon)^r$ . The proof of this fact is given

in [FR3] and the main tool in proving it is a rather difficult theorem about intersection properties of partitions [FR1].

Concerning the original problem of characterizing Ramsey configurations, the first problem left in [FR3] was to decide whether the vertices of the regular pentagon form a Ramsey configuration. Recently, I. Kriz [K] proved that the vertices of the regular  $k-gon$   $A_k$  form a Ramsey configuration for any positive integer  $k$ . This supports the hope that Ramsey configurations might coincide with spherical ones. The bounds on  $n(A_k, r)$  which follow from the proof given in [K], seem to be rather weak.

Finally, we will state here one more problem. In [E5], P. Erdős asked whether  $n(T, 2) = 2$  for any configuration of three vertices of a non-regular triangle.

## Van der Waerden's Theorem

In 1927, B.L. van der Waerden [V] published the proof of the following result:

**Theorem (Van der Waerden's Theorem).** *For all positive integers  $k$  and  $r$ , there exists an integer  $W(k, r)$  so that if the set of integers  $\{1, 2, \dots, W(k, r)\}$  is partitioned into  $r$  classes, then at least one class contains a  $k$ -term arithmetic progression.*

Set  $W(k) = W(k, 2)$ , known values are  $W(2) = 3$ ,  $W(3) = 9$ ,  $W(4) = 35$ ,  $W(5) = 178$  (see [Gr2]). The original proof of van der Waerden, as well as the later proof (cf. [GR3], [D3], [T1]), give very weak upper bounds on the function  $W(k)$ . In 1987, S. Shelah [S2] found a fundamentally new proof that yields a much better upper bound.

Set  $f_1(n) = 2n$ ,  $f_2(n) = 2^n$ ,  $\dots$ ,  $f_{i+1}(n) = f_i^n(1)$ , where  $f_i^n$  denotes the  $n$ -times iterated function  $f_i$ . So,  $f_3(n)$  is the  $n$ -times iterated exponential consisting of a tower of  $n$  twos. The function  $f_4(n)$  grows much faster:  $f_4(1) = f_3(1) = 2$ ,  $f_4(2) = f_3 \circ f_3(1) = 2^2$ ,  $f_4(3) = f_3 \circ [f_3 \circ f_3(1)] = f_3[f_4(2)] = 65536$  which is a tower of 4 twos, and  $f_4(4) = f_3[f_4(3)]$  is a tower of 65536 twos. While the earlier best upper bounds on  $W(k)$  were of the order of the Ackerman function, i.e.  $f_k(ck)$  for  $c > 0$ , the bound given by S. Shelah [S2] is "only"  $W(k) \leq f_4(k + 1)$  (cf. also [GRS]). The problem of finding reasonable estimates for  $W(k)$  is one of the main problems of the area and the result of Shelah is a vast improvement of the previous bound. The best known lower bounds are exponential:  $W(k + 1) > k \cdot 2^k$  for  $k$  prime was proved in [B2], while  $W(k) > k^{1-\epsilon} 2^k$  for  $k > k_0(\epsilon)$  was recently established in [S6].

In [ET], P. Erdős and P. Turán considered the quantity  $r_k(n)$  defined as the maximum cardinality  $|Z|$  of a set  $Z \subset \{1, 2, \dots, n\}$  that does not contain an arithmetic progression of  $k$  terms. One of the reasons for investigating  $r_k(n)$  is that  $r_k(n) < n/2$  implies  $W(k) < n$  and bounds on  $r_k(n)$  could improve the upper bound for  $W(k)$ . Erdős and Turán noted that  $r_k(n + m) \leq r_k(n) + r_k(m)$  which implies that  $\lim_{n \rightarrow \infty} c_k/n = c_k$  exists. Erdős and Turán conjectured that  $c_k = 0$  for all  $k$ . They also conjectured that  $r_k(n) < n^{1-\epsilon_k}$  which was disproved by Salem and Spencer [SS] who proved  $r_3(n) > n \exp\left(-\frac{c \ln n}{\ln \ln n}\right)$ . This was improved in 1946 by F.A. Behrend who

established  $r_3(n) > n \exp\left(-\frac{c \ln n}{\sqrt{\ln n}}\right)$ , the best current upper bound for  $r_3(n)$  is due to Heath-Brown [H2] and Szemerédi. In 1952, Roth proved that  $r_3(n) < c n/\log \log n$ . In 1967, E. Szemerédi [S7] proved  $r_4(n) = o(n)$ . His proof used the general statement of van der Waerden's theorem. Later, Roth [R8, R9] gave a different proof that  $r_4(n) = o(n)$  (and his proof probably yields  $r_4(n) < n/\log_l n$  where  $l$  is a large fixed integer and  $\log_l n$  denotes the  $l$ -times iterated logarithm cf [S8]). In 1973, Szemerédi [S9] gave a proof of the Erdős-Turán conjecture, establishing that  $c_k = 0$  for every  $k$ . A different proof based on ergodic theory was found in 1977 by H. Fürstenberg [Fu]. Since then, the methods of ergodic theory have proven to be a powerful tool for investigating related questions; proofs of several extensions of Szemerédi's theorem were found [FK1]. Recently, Fürstenberg and Katzenelson, establishing a conjecture of R.L. Graham [Gr1], proved the density form of Hales-Jewett's theorem. The methods of ergodic theory as well as the proof of Szemerédi do not yield bounds on  $r_k(n)$ . Therefore, finding alternative proof techniques is still of interest. Recently, P. Frankl and the author considered the following question. Let  $\{a_i, b_i\}, 1 \leq i \leq k$  be pairwise disjoint 2-element sets. Define  $F_i = \{a_1, a_2, \dots, a_k, b_i\} - \{a_i\}$  and  $F = F(k) = \{F_1, F_2, \dots, F_k\}$ . Let  $\mathcal{E} \subset [V]^k$  be such that i)  $|V| = n$ , ii)  $|E_1 \cap E_2| \leq k - 2$ , for any two distinct  $E_1, E_2 \in \mathcal{E}$  and iii)  $\mathcal{E}$  does not contain  $F(k)$  as a subconfiguration.

Set  $e\tilde{x}(n, F(k)) = \text{Max}\{|\mathcal{E}|; \mathcal{E} \text{ satisfies i), ii), and iii)}\}$ . With P. Frankl, we recently noticed that  $e\tilde{x}(n, F(k)) \geq c_k n^{k-2} r_k(n)$ , and thus,

$$e\tilde{x}(n, F(k)) = o(n^{k-1}) \quad (6)$$

would imply  $r_k(n) = o(n)$ . (In fact, this implies a density version of the Graham-Leeb-Rothschild Theorem as well.) This is known to be true for  $k = 3$  by [RS] and with P. Frankl, we recently gave the proof for  $k = 4$ . We conjecture that (6) holds for every  $k$ .

## References

- [AH] F.G. Abramson, L.A. Harrington: Models without indiscernibles. *J. Symbolic Logic* **43** (1978) 572–600
- [AKS] M. Ajtai, J. Komlós, E. Szemerédi: A note on Ramsey numbers. *J. Comb. Th. (A)* **29** (1980) 354–360
- [B1] F.A. Behrend: On sets of integers which contain no three in arithmetic progression. *Proc. Nat. Acad. Sci.* **33** (1946) 331–332
- [B2] E.R. Berlekamp: A Construction for partitions which avoid long arithmetic progressions. *Can. Math. Bull.* **11** (1968) 409–414
- [CS] T.J. Carlson, S.G. Simpson: Topological Ramsey theory. To appear in: *Mathematics of Ramsey theory* (J. Nešetřil and V. Rödl, eds.). Springer, Berlin
- [CW] D. Coppersmith, S. Winograd: Matrix multiplication via arithmetic progressions. *19th Annual ACM Symposium on Theory of* (1987)
- [D1] W. Deuber: Generalization of Ramsey's theorem. In: *Infinite and Finite Sets* (A. Hajnal, R. Rado and V.T. Sós, eds.) *Colloq. Math. Soc. János Bolyai* **10**. North-Holland, Amsterdam 1975, pp. 323–332
- [D2] W. Deuber: Partitionen und lineare Gleichungssysteme. *Math. Z.* **133** (1973) 109–123

- [D3] W. Deuber: On van der Waerden's theorem on arithmetic progressions. *J. Comb. Th.* **A32** (1982) 115–118
- [DR] W. Deuber, B.L. Rothschild: Categories without the Ramsey property. In: *Combinatorics* (A. Hajnal, V.T. Sós, eds.) *Coll. Math. Soc. János Bolyai*. **18** (1976) 225–249
- [E1] P. Erdős: Some remarks on the theory of graphs. *Bull. Amer. Math. Soc.* **53** (1947) 292–294
- [E2] P. Erdős: Graph theory and probability. *Canad. J. Math.* **11** (1959) 34–38
- [E3] P. Erdős: The art of counting. *Selected Writings*. The MIT Press, Cambridge, Massachusetts and London, England (1973)
- [E4] P. Erdős: Problems and results on finite and infinite graphs. In: *Recent advances in graph theory*. Czechoslovak Academy of Sciences 1975, pp. 183–192
- [E5] P. Erdős: Combinatorial problems in geometry and number theory. In: *Proceedings of the Symposium in Pure Math*, vol. 34 (1979) AMS, pp. 149–162
- [EG] P. Erdős, R.L. Graham, P. Montgomery, B.L. Rothschild, J. Spencer, E.G. Straus: Euclidean Ramsey Theorems. I. *J. Combin. Theory Ser. A* **14** (1973) 341–363
- [EHP] P. Erdős, A. Hajnal, L. Pósa: Strong embeddings of graphs into colored graphs. In: *Infinite and finite sets* (A. Hajnal, R. Rado and V.T. Sós, eds.) *Colloq. Math. Soc. János Bolyai* 10. North-Holland, Amsterdam 1975, pp. 585–595
- [EHMR] P. Erdős, P. Hajnal, A. Máté, R. Rado: Combinatorial set theory. Partition relations for cardinals. North-Holland, Amsterdam 1984
- [EHR] P. Erdős, A. Hajnal, R. Rado: Partition relations for cardinal numbers. *Acta Math. Acad. Sci. Hungar.* **16** (1965) 93–196
- [EL] P. Erdős, L. Lovász: Problems and results on 3-chromatic hypergraphs and some related questions. In: *Infinite and finite sets* (A. Hajnal, R. Rado and V.T. Sós, eds.), North-Holland, Amsterdam 1975, pp. 609–628
- [ES] P. Erdős, G. Szekeres: A Combinatorial problem in geometry. *Comp. Math.* **2** (1935) 463–470
- [ET] P. Erdős, and P. Turán: On some sequences of integers. *J. London Math. Soc.* **11** (1936) 261–264
- [Fa] K.J. Falconer: The realization of distances in measurable subsets covering  $R^n$ . *J. Combin. Theory, Ser A* **31** (1981) 184–189
- [Fo] J. Folkman: Graph with monochromatic complete subgraphs in every edge coloring. *SIAM J. Appl. Math.* **18** (1970) 19–24
- [FGR] P. Frankl, R.L. Graham, V. Rödl: Induced restricted Ramsey theorems for spaces. *J. Comb. Theory, Ser. A* **44** (1987) 120–128
- [FR1] P. Frankl, V. Rödl: Forbidden Intersections. *Trans. Amer. Math. Soc.* **299** (1987)
- [FR2] P. Frankl, V. Rödl: All triangles are Ramsey. *Trans. Amer. Math.* **297** (1986) 777–779
- [FR3] P. Frankl, V. Rödl: A partition property of simplices in Euclidean space. *J. Amer. Math.* **3**(1) (1990) 1–7
- [FR4] P. Frankl, V. Rödl: The uniformity lemma for hypergraphs. *Graphs and Combinatorics*. (To appear)
- [FW] P. Frankl, R.M. Wilson: Intersection theorems with geometric consequences. *Combinatorica* **1** (1981) 357–368
- [Fr] Frasnay: Sur des fonctions d'entiers se rapportant au Ramsey. *C. R. Acad. Sci. Paris* **256** (1963) 2507–2510
- [Fu] H. Fürstenberg: Ergodic behavior of diagonal measures and a theorem of Szemerédi on arithmetic progressions. *J. Anal. Math.* **31** (1977) 204–256
- [FK1] H. Fürstenberg, Y. Katznelson: An ergodic Szemerédi theorem. IP-Systems and Combinatorial Theory. *J. Anal. Math.* **45** (1985) 117–168
- [FK2] H. Fürstenberg, Y. Katznelson: A density version of the Hales-Jewett theorem. (1990)

- [FW] P. Frankl, R.M. Wilson: Intersection theorems with geometric consequences. *Combinatorica* **1** (1981) 357–368
- [Go] Goodman: On sets of acquaintances and strangers at any party. *Amer. Math. Monthly* **66** (1959) 778–783
- [Gr1] R.L. Graham: Old and new euclidean Ramsey theorems. Discrete geometry and convexity. *Ann. New York Acad. Sci.* **440** (1985) 20–30
- [Gr2] R.L. Graham: Recent developments in Ramsey theory. *Proceedings of the International Congress of Mathematicians, Warszawa 1983*, pp. 1555–1567
- [GLR] R.L. Graham, K. Leeb, B.L. Rothschild: Ramsey's theorem for a class of categories. *Adv. Math.* **8** (1972) 417–433
- [GR1] R.L. Graham, V. Rödl: Numbers in Ramsey theory. In: *Surveys in Combinatorics* (ed. C. Whitehead). (LMS Lecture Note Series 123.) Cambridge University Press, pp. 111–153
- [GR2] R.L. Graham, B.L. Rothschild: Ramsey's theorem for  $n$ -parameter sets. *Trans. Amer. Math. Soc.* **159** (1971) 257–292
- [GR3] R.L. Graham, B.L. Rothschild: A short proof of van der Waerden's theorem on arithmetic progressions. *Proc. Amer. Math. Soc.* **42** (1974) 356–386
- [GRS] R.L. Graham, B.L. Rothschild, J.H. Spencer: *Ramsey Theory*. John Wiley and Sons, New York 1980
- [GY] J.E. Graver, J. Yackel: Some graph theoretic results associated with Ramsey's theorem. *J. Comb. Theory Ser. A* **4** (1968) 125–175
- [H1] H. Hadwiger: Ungelöste Probleme. No. 40. *Elemente der Math.* **16** (1961) 103–104
- [H2] D.R. Heath-Brown: Integer sets containing no arithmetic progressions. 1986
- [HJ] A.W. Hales, R.I. Jewett: Regularity and positional games. *Trans. Amer. Math. Soc.* **106** (1963) 222–229
- [K] I. Kriz: Permutation groups in Euclidean Ramsey theory. (To appear)
- [KV] P.G. Kolaitis, M. Vardi: The decision problem for the probabilities of higher-order properties. In: 19th 'Annual ACM Symposium (STOC)', pp. 425–435
- [L1] K. Leeb: Vorlesungen über Pascaltheorie. Erlangen 1973
- [L2] G. Lorden: Blue-empty chromatic graphs. *Amer. Math. Monthly* **69** (1962) 114–120
- [L3] L. Lovász: Self-dual polytopes and the chromatic number of distance graphs on the sphere. *Acta Sci. Math.* **45** (1983) 317–323
- [LR] D.G. Larman, C.A. Rogers: The realization of distances within sets in Euclidean space. *Mathematika* **19** (1972) 1–24
- [MM] L. Moser, W. Moser: Solution to Problem 10. *Cand. Math. Bull.* **16** (1961) 187–189
- [N1] J. Nesetril: There are 4 types of Ramsey classes of graphs. *J. Comb. Theorems B* **46**(2) (1989) 127–132
- [N2] J. Nesetril: Ramsey theory. (To appear in *Handbook of Combinatorics*, R.L. Graham, M. Grötschel, L. Lovász, ed.)
- [NR1] J. Nesetril, V. Rödl: Type theory of partition properties of graphs. In: *Recent Advances in Graph Theory* (ed. M. Fiedler). Academia, Prague 1975, pp. 405–412
- [NR2] J. Nesetril, V. Rödl: A structural generalization of the Ramsey theorem. *Bull. Amer. Math. Soc.* **83**(1) (1977) 127–128
- [NR3] J. Nesetril, V. Rödl: Partition of relational and set-systems. *J. Comb. Theorems A* **22** (1977) 289–312
- [NR4] J. Nesetril, V. Rödl: A short proof of the existence of restricted Ramsey graphs by means of a partite construction. *Combinatorica* **1**(2) (1981) 199–202
- [NR5] J. Nesetril, V. Rödl: Ramsey classes of set systems. *J. Comb. Theory A* **34** (1983) 183–201
- [NR6] J. Nesetril, V. Rödl: Combinatorial partitions of finite cosets and Lattices – Ramsey lattices. *Algebra Universalis* **19** (1984) 106–119

- [NR7] J. Nešetřil, V. Rödl: Sparse Ramsey graphs. *Combinatorica* **4**(1) (1984) 71–78
- [NR8] J. Nešetřil, V. Rödl: Strong Ramsey theorems for Steiner systems. *Trans. Amer. Math. Soc.* **303**(1) (1987) 183–197
- [NR9] J. Nešetřil, V. Rödl: Partite construction and Ramseyian theorems for sets, numbers and spaces. *Comment. Math. Univ. Carol.* **28** (1987) 569–580
- [NR10] J. Nešetřil, V. Rödl: The partite construction and Ramsey set systems. *Discr. Math.* **75** (1989) 327–334
- [NR11] Mathematics of Ramsey theory (ed. J. Nešetřil and V. Rödl). (Algorithms and Combinatorics, vol. 5). Springer, Berlin Heidelberg New York 1990
- [NR12] J. Nešetřil, V. Rödl: The partite construction and Ramsey space systems. In: Mathematics of Ramsey theory. Springer, Berlin Heidelberg New York 1990
- [P1] H. J. Prömel: Induced partition properties of combinatorial cubes. *J. Comb. Theory Ser. A* **39** (1985) 177–208
- [P2] H. J. Prömel: Partition properties of  $q$ -hypergraphs. *J. Comb. Theory Ser. B* **41** (1986) 356–385
- [PH] J. Paris, L. Harrington: A mathematical incompleteness in Peano arithmetic. In: *Handbook of Mathematical Logic* (ed. J. Barwise). North-Holland, Amsterdam 1977, pp. 1133–1142
- [PV] H.J. Prömel, B. Voigt: A sparse Graham-Rothschild theorem. *Trans. Amer. Math. Soc.* **309**(1) (1988) 113–137
- [RK] S.P. Radziszowski, D.L. Kreher: Search algorithm for Ramsey graphs by union of group orbits. *J. Graph Theory* **12**(1) (1988) 59–72
- [R1] R. Rado: Studien zur Kombinatorik. *Math. Z.* **36** (1933) 425–480
- [R2] R. Rado: Note on combinatorial analysis. *Proc. London Math. Soc.* **48** (1943) 122–160
- [R3] F.P. Ramsey: On a Problem for Formal Logic. *Proc. London Math. Soc.* **30** (1930) 264–286
- [R4] V. Rödl: A generalization of the Ramsey theorem. In: *Graphs, hypergraphs and block systems*, Zielona Góra 1976, pp. 211–219
- [R5] V. Rödl: On a problem in combinatorial geometry. *Discr. Math.* **45** (1983) 129–131
- [R6] V. Rödl: Upper bounds on Ramsey numbers  $R(k, l)$ . To appear
- [R7] K. Roth: On Certain Sets of Integers. *J. London Math.* **28** (1953) 104–109
- [R8] K. Roth: Irregularities of Sequences Relative to Arithmetic Progressions, III. *J. Number Theory* **41** (1970) 125–142
- [R9] K. Roth: Irregularities of sequences relative to arithmetic progressions, IV. *Per. Math. Hungar.* (1972) 301–326
- [R10] S.P. Radziszowski: Small Ramsey numbers (manuscript 1991)
- [RS] I.Z. Ruzsa, E. Szemerédi: Triple systems with no six points carrying three triangles. *Colloq. Math. J. Bolyai* **18** (1978) 939–945
- [SS] R. Salem, D.C. Spencer: On sets of integers which contain no three terms in arithmetic progression. *Proc. Nat. Acad. Sci.* **28** (1942) 561–563
- [S1] I. Schur: Über die Kongruenz  $x^m + y^m \equiv z^m \pmod{\varrho}$ . *Jber. Deutsch. Math.-Verein.* **25** (1916) 114–117
- [S2] J.B. Shearer: A note on the independence number of a triangle-free graph. *Discr. Math.* **46** (1983) 83–87
- [S3] S. Shelah: Primitive recursive bounds for van der Waerden numbers. *J. Amer. Math. Soc.* **1** (1988) 683–697
- [S4] J. Spencer: Ramsey's theorem – A new lower bound. *J. Comb. Theory (A)* **18** (1975) 108–115
- [S5] J. Spencer: Asymptotic lower bounds for Ramsey functions. *Discr. Math.* **20** (1977) 69–76

- [S6] Z. Szabó: An application of Lovász local lemma – A new lower bound for the van der Waerden number. *Random structures and algorithms* **1**(3) (1990) 343–360
- [S7] E. Szemerédi: On sets of integers containing no four elements in arithmetic progression. *Acta Math. Acad. Sci. Hung.* **20** (1969) 89–104
- [S8] E. Szemerédi: On sets of integers containing no  $k$ -elements in arithmetic progression. *Proceedings of the International Congress of Mathematicians, Vancouver 1974*, pp. 503–505
- [S9] E. Szemerédi: On sets of integers containing no  $k$  elements in arithmetic progression. *Acta Arith.* **27** (1975) 199–245
- [T1] A.D. Taylor: A note on van der Waerden theorem. *J. Comb. Theory Ser. A* **33** (1982) 215–219
- [T2] A. Thomasson. Upper bounds for Ramsey numbers (to appear)
- [V] B.L. van der Waerden: Beweis einer Baudetschen Vermutung. *Nieuw Arch. Wisk.* **15** (1927) 212–216
- [W] D.R. Woodal: Distances realized by sets covering plane. *J. Comb. Theory A* **14** (1973) 187–200

# Strongly Polynomial and Combinatorial Algorithms in Optimization

Éva Tardos

School of Operations Research and Industrial Engineering  
Cornell University, Ithaca, NY 14853, USA

## 1. Introduction

The input to a combinatorial optimization problem usually consists of two parts: the first part describes the combinatorial structure; the second part is a list of numerical data. For example, the input to the maximum-flow and the shortest path problems consists of a network (the combinatorial structure) and numbers that define the capacity and the length of each arc, respectively. An algorithm for these problems is *polynomial* if its running time can be bounded by a polynomial in the size of the underlying combinatorial structure (the number of nodes and edges) and the number of digits needed to write the numerical data. Most of the early polynomial algorithm satisfied, in fact, a stronger notion of efficiency; these algorithms are not only polynomial, but the number of arithmetic operations performed can be bounded by a polynomial in the size of the underlying combinatorial structure alone, independent of the size of the numbers involved. Such algorithms are called *strongly polynomial*.

Strongly polynomial algorithms are more appealing theoretically and the additional insight needed to develop such algorithms can lead to practical improvements. It is an important theoretical question to understand which problems can be solved in strongly polynomial time. In this paper we shall survey partial results in this direction. We discuss some techniques to turn polynomial algorithms into strongly polynomial ones.

Khachiyan's [15] proof that the linear programming problem,  $\max(cx : Ax \leq b)$ , can be solved in polynomial time has been one of the major breakthroughs in the design and analysis of algorithms in the last ten years. The first polynomial linear programming algorithm was developed using the ellipsoid method of Yudin and Nemirovskii [33] from non-linear programming. Grötschel, Lovász and Schrijver [11] have used the ellipsoid method to develop a powerful technique for proving that many combinatorial optimization problems can be solved in polynomial time (see also [14] and [25]). The polynomial algorithms developed using the ellipsoid method are generally not strongly polynomial. There have been several other polynomial linear programming algorithms developed in the recent years (e.g., Karmarkar's algorithm [13]), however no strongly polynomial

is known. It is an important open problem if there exists a strongly polynomial linear programming algorithm.

A somewhat related question is if the combinatorial optimization problems that were proved to be solvable in polynomial time using the ellipsoid method can be solved more efficiently using combinatorial techniques, rather than employing the analytic techniques from non-linear programming.

In this paper we shall survey two techniques for converting polynomial algorithms into strongly polynomial ones. The first one uses Diophantine approximation. It replaces the numbers occurring in the problem description by small integers that define an equivalent problem. Any polynomial algorithm can be used to solve the resulting problem in strongly polynomial time. The second technique is iterative. After some preprocessing each iteration finds an approximately optimal solution to the problem. From such an approximate solution the algorithm either finds an optimal solution, or it concludes that some of the defining constraint are not necessary, and these are deleted before the next iteration starts.

There are further techniques known for converting some polynomial algorithms into strongly polynomial ones. The first such technique has been developed by Megiddo [20] for solving combinatorial ratio minimization problems in strongly polynomial time. This technique can turn algorithms that use binary search (which is not strongly polynomial) into strongly polynomial algorithms. One of the most powerful application of this technique, due to Megiddo [21], is testing the feasibility of linear programs with at most two variables in each inequality. A more recent application by Norton, Plotkin and Tardos [22] gives a further extension of the class of linear programs solvable in strongly polynomial time: If a linear program is known to be solvable in strongly polynomial time, then so is its extension by a constant number of additional variables and side constraints. A nice example in this class, which we shall discuss later in more detail, is the concurrent flow problem.

## 2. Using Diophantine Approximation

One of the most powerful techniques for making combinatorial optimization algorithms strongly polynomial, due to Frank and Tardos [7], uses simultaneous Diophantine approximation. Many combinatorial optimization problems can be defined as a pair of a (highly structured) system of subsets  $\mathcal{I}$  of a finite set  $E$  and weights  $w(e)$  for each  $e \in E$  (e.g., perfect matchings in a graph form a system of the subsets of the edges). The weight of a set  $I \in \mathcal{I}$  is  $w(I) = \sum_{e \in I} w(e)$ . The problem is to find a set in  $\mathcal{I}$  with maximum weight. Let  $n$  denote the size of the set  $E$ . The size of such the combinatorial structure  $\mathcal{I}$  depends on the way  $\mathcal{I}$  is specified, but generally it is at least  $n$ .

The idea of converting polynomial algorithms for problems in the above form into strongly polynomial ones is as follows. Consider first a weight function  $\hat{w}$  whose coordinates are integers with size no more than a polynomial in  $n$ . For such inputs there is no distinction between polynomial and strongly polynomial

algorithms. The key lemma is that for every weight function  $w$  there is an equivalent weight function  $\hat{w}$  that consists of small integers. Furthermore, the new weight function  $\hat{w}$  can be found in strongly polynomial time. The strongly polynomial algorithm first computes such an equivalent weight function  $\hat{w}$ , and then uses the polynomial algorithm to find the maximum weight set subject to the new weights.

Two weight functions are *equivalent* for a problem  $\mathcal{I}$ , if for every two sets  $I, J \in \mathcal{I}$  we have that  $w(I) < w(J)$  if and only if  $\hat{w}(I) < \hat{w}(J)$ . Clearly, the maximum weight set in  $\mathcal{I}$  is the same for any two equivalent weight functions. The existence of an equivalent weight function whose coordinates are small integers was first observed by Orlin [23]. Orlin's proof is not algorithmic. Frank and Tardos [7] gave an algorithmic proof using the simultaneous Diophantine approximation technique of Lovász (see in [18]).

Throughout the paper we shall use different norms. For a vector  $x$  we use  $\|x\|_\infty$  to denote the maximum absolute value of a coordinate of  $x$ , and  $\|x\|_1$  to denote the sum of the absolute values of the coordinates of  $x$ .

**Theorem 1.** *For a given  $n$ -dimensional vector  $w$ , and an integer  $N$  one can find an integer vector  $\hat{w}$  in time polynomial in  $n$  and  $\log N$ , such that  $\|\hat{w}\|_\infty \leq 2^{O(n^3)}N^{O(n^2)}$  and for every integer vector  $a$  such that  $\|a\|_1 \leq N$ ,  $aw < 0$  if and only if  $a\hat{w} < 0$ .*

Let  $w$  be the weight function in a combinatorial optimization problem. If Theorem 1 is applied with  $w$  and  $N = n + 1$ , we obtain an equivalent weight function  $\hat{w}$ .

Next we try to give an idea of the proof of Theorem 1. Consider a vector  $w$  and an integer  $N$ . Lovász's [18] algorithm finds a positive integer  $q$  and an integer vector  $w'$  such that (1)  $\|gw - w'\|_\infty \leq 1/N$ , and (2)  $q \leq 2^{O(n^2)}N^n$ . This vector  $w'$  satisfies a property similar to that required by Theorem 1.

**Lemma 2.** *Let  $w'$  be the vector obtained from  $w$  and  $N$  by the Diophantine approximation algorithm, and let  $a$  be an integer vector such that  $\|a\|_1 \leq N$ . If  $aw \leq 0$  then also  $aw' \leq 0$ .*

*Proof.* Assume  $aw \leq 0$ . Consider  $aw'$ . The sum of the coordinates of  $a$  is at most  $N$ . Using this fact, and the first property of Diophantine approximation we obtain that  $a(w' - qw) < 1$ . We also have that  $aw \leq 0$  and  $q \geq 0$ , therefore  $aw' < 1$ . Since both  $a$  and  $w'$  are integral, it follows that  $aw' \leq 0$ .  $\square$

The lemma implies that the vector  $w'$  satisfies all the relevant non-strict inequalities satisfied by  $w$ . The vector  $\hat{w}$  in Theorem 1 must also satisfy all of the strict inequalities. Theorem 1 can be proved by a procedure that uses Diophantine approximation repeatedly.

A nice application of Theorem 1 shows that a maximum weight clique in a perfect graph can be found in strongly polynomial time. The first polynomial algorithm for the problem was developed by Grötschel, Lovász and Schrijver [11] using the ellipsoid method. Theorem 1 can be used to convert this algorithm into a strongly polynomial one.

Many combinatorial optimization problems are related to linear programming. Combinatorial special cases of linear programming include several network flow problems. These problems generally have constraint matrix with small integer entries. In Section 4 we shall discuss the maximum-flow, the transshipment and the multi-commodity flow problems in more detail. The constraint matrix in these problems describes the combinatorial structure (e.g., it is the incidence matrix of the network). The simultaneous Diophantine approximation technique can be used to convert polynomial algorithms for these problems into strongly polynomial ones.

Consider the linear program  $\max(cx : Ax \leq b)$  where  $A$  is an  $n$  by  $m$  matrix. Assume that the matrix  $A$  is integral,  $m \geq n$ , and let  $\Delta(A)$  denote the maximum absolute value of a subdeterminant of  $A$ . We say that two objective functions  $c$  and  $\hat{c}$ , and the right-hand sides  $b$  and  $\hat{b}$  are equivalent for the constraint matrix  $A$  if the same set of inequalities are satisfied as equations at the optimal solution for the two linear programs  $\max(cx : Ax \leq b)$  and  $\max(\hat{c}x : Ax \leq \hat{b})$ . Using linear programming duality and Theorem 1 we can prove the following theorem.

**Theorem 3.** *For any linear program  $\max(cx : Ax \leq b)$  with an  $n$  by  $m$  integer matrix  $A$ , there is an equivalent integral right-hand side  $\hat{b}$  and objective function  $\hat{c}$  such that  $\log \|b\|_\infty$  and  $\log \|c\|_\infty$  are polynomially bounded in  $n$ ,  $m$  and  $\log \Delta(A)$ .*

**Corollary 4.** *Linear programs where the constraint matrix  $A$  is integral and has coefficients whose size is at most polynomial in  $n$  and  $m$ , can be solved in strongly polynomial time.*

### 3. Iterative Approach

In this section we discuss an iterative method for converting polynomial algorithms into strongly polynomial ones. When using the iterative technique it is possible to take advantage of additional insight into the combinatorial structure of the problem in question. In some cases this leads to very efficient algorithms. For example, highly efficient algorithms for the maximum-flow and transshipment problems have been obtained this way.

Consider a linear program  $\max(cx : Ax \leq b)$ . The *dual* of this problem is the problem  $\min(by : A^T y = c, y \geq 0)$ . We say that  $x'$  is *feasible* if  $Ax' \leq b$ , and  $y'$  is *dual feasible* if  $y' \geq 0$  and  $A^T y' = c$ . It is well-known that the linear program and its dual have the same optimal value. *Complementary slackness* gives a localized condition that helps recognize a pair of primal and dual optimal solutions. A feasible solution  $x'$  and a feasible dual solution  $y'$  are optimal if and only if for every row  $a_i x' \leq \beta_i$  of  $Ax \leq b$  and the corresponding coordinate  $y_i$  of  $y$ , we have

$$a_i x' < \beta_i \implies y'_i = 0 . \quad (1)$$

The key concept of the iterative method is  $\varepsilon$ -optimality. This is a relaxation of the complementary slackness conditions. A feasible solution  $x'$  and feasible dual solution  $y'$  are  $\varepsilon$ -optimal if for every row  $i$  we have that

$$a_i x' < \beta_i - \varepsilon \implies y'_i = 0 . \quad (2)$$

Note that  $\varepsilon$ -optimality with  $\varepsilon = 0$  is the same as the complementarity slackness conditions. It is not too difficult to show that an  $\varepsilon$ -optimal pair  $(x', y')$  is close to being optimal both in terms of its objective function value and in terms of the distance to an optimal solution. This was first proved by Mangasarian [19]. For our purposes it will be important that the bound on the distance of the  $\varepsilon$ -optimal  $x'$  to an optimal solution is independent of the size of the numbers in the objective function. Such a bound was proven by Cook, Gerards, Schrijver and Tardos [4].

**Theorem 5.** *Suppose  $x'$  and  $y'$  form a pair of  $\varepsilon$ -optimal primal and dual solutions to the linear program  $\max(cx : Ax \leq b)$ . Then there exists an optimal  $x^*$  such that  $\|x' - x^*\|_\infty \leq n\varepsilon\Delta(A)$ .*

Assume that  $A$  is an integer matrix, and let  $\alpha$  denote  $\|A\|_\infty$ , the largest absolute value of a coefficient of  $A$ .

**Corollary 6.** *Suppose that  $x'$  and  $y'$  are a pair of  $\varepsilon$ -optimal solutions to the linear program  $\max(cx : Ax \leq b)$  and its dual. Consider a constraint  $ax \leq \beta$ . Suppose  $x'$  satisfies  $\beta - ax' \geq n\varepsilon\Delta(A)$ . The coefficient corresponding to the inequality  $ax \leq \beta$  is zero in every optimal dual solution, and every optimal solution is also optimal if the constraint  $ax \leq \beta$  is deleted.*

This theorem and corollary are the key of the iterative strongly polynomial algorithm. The algorithm repeatedly finds an  $\varepsilon$ -optimal primal-dual solution pair, and deletes all constraints that are known (by Corollary 6) not to be tight at an optimal solution. This procedure is repeated until only the tight constraints are left. The algorithm terminates in at most as many iterations as the number of defining inequalities. There are two remaining issues that need to be discussed: how can one guarantee that at least one inequality will be deleted each iteration, and how can one find an  $\varepsilon$ -optimal solution in strongly polynomial time.

In order to guarantee that at least one inequality will be deleted one has to do some preprocessing before the iteration, and choose an appropriate  $\varepsilon$ . In most special cases the preprocessing is quite simple, in the general case it involves a projection. See [30, 27] for more details for the general case. After such preprocessing the appropriate choice of  $\varepsilon$  turns out to be  $\|b\|_\infty/(n^2\Delta(A))$ .

Next consider the issue of finding an  $\varepsilon$ -optimal solution. This can be accomplished in strongly polynomial time if  $\Delta(A)$  is “small”. Round every coordinate of  $b$  to an integer multiple of  $\varepsilon$ . A pair of optimal primal and dual solutions for the rounded problem is  $\varepsilon$ -optimal for the original problem. If  $\Delta(A)$  is small then the coefficients of the rounded  $b$  are small multiples of  $\varepsilon$ . If a polynomial, but not strongly polynomial, algorithm is used to solve the rounded problem, then the resulting running time will be independent of the size of the numbers in the original vector  $b$ . However, it will depend on the size of the numbers in the objective function,  $c$ . In order to get rid of the dependence on  $c$  we need to solve the rounded problem by applying the same iterative scheme to its dual. The running time of the resulting algorithm is independent of the size of the numbers in the objective function as well as in the right-hand side.

## 4. Network Flow Problems

The most well-known linear programs with small integral constraint matrices are network flow problems. These problems can be used to illustrate some of the ideas from the previous sections. By taking advantage of the special structure one can make the general algorithms simpler and more efficient.

A network is a directed graph  $G = (V, E)$ . We shall use  $n$  and  $m$  to denote the number of nodes and edges, respectively. The *maximum-flow problem* is defined by a network, a source  $s \in V$  and a sink  $t \in V$  and nonnegative capacities  $u(v, w)$  associated with the edges  $(v, w) \in E$ . A vector  $f(v, w)$  for  $(v, w) \in E$  is a *preflow* if  $0 \leq f(v, w) \leq u(v, w)$  for every  $e \in E$ . Its *excess* at a node  $v \in V$  is defined to be

$$e_f(v) = \sum_w f(w, v) - \sum_w f(v, w) . \quad (3)$$

A preflow  $f$  is a *flow* if  $e_f(v) = 0$  for every node  $v \neq s, t$ . The *value* of the flow is  $e_f(t)$ . The problem is to find a flow with maximum value.

In the *transshipment problem* there are costs  $c(v, w)$  associated with the edges instead of capacities, and there are demands  $b(v)$  on the nodes  $v \in V$ . A preflow  $f$  is a *transshipment* if  $e_f(v) = b(v)$  for every node  $v \in V$ . The *cost of a preflow*  $f$  is  $\sum_{vw \in E} f(v, w)c(v, w)$ . The problem is to find a transshipment of minimum cost.

Both the maximum-flow and the transshipment problems are linear programs with  $0, \pm 1$  constraint matrix. The capacities and the demands form the right-hand side, the costs are the coefficients of the objective function. The strongly polynomial solvability of these problems follows, for example, from Corollary 4.

The first polynomial algorithm for the maximum-flow problem is due independently to Dinic [5] and Edmonds and Karp [6]. Both algorithms are strongly polynomial. The first polynomial algorithm for the transshipment problem is due to Edmonds and Karp [6]. This algorithm is *not* strongly polynomial. The first strongly polynomial algorithm was developed much later [29] using the iterative method of Section 3.

The *multi-commodity flow problem* is defined similarly. There are  $k$  pairs of sources and sinks  $s_i, t_i \in V$  and each has an associated demand,  $d_i$ . The problem is to find flows  $f_i$  from  $s_i$  to  $t_i$  of value  $d_i$  such that a joint capacity constraint,  $\sum_{(v,w) \in E} f_i(v, w) \leq u(v, w)$ , is satisfied for each edge  $(v, w) \in E$ . The multi-commodity flow problem is also a linear program with  $0, \pm 1$  constraint matrix, and therefore it can be solved in strongly polynomial time.

The multi-commodity flow problem is fairly a general linear program. Dinic (see in [1]) and Itai [12] have independently proved that any linear program can be reduced to a 2-commodity flow problem in polynomial time.

### 4.1 Transshipment Problem

The first strongly polynomial algorithm for the transshipment problem [29] is rather slow and complicated. It is more appropriate to refer to it as a proof of solvability in strongly polynomial time, rather than as an efficient algorithm. Since then, many strongly polynomial algorithms have been discovered and some

of them are surprisingly simple and efficient, such as the recent algorithms due to Goldberg and Tarjan [9] and Orlin [24]. Here we shall focus Orlin's algorithm.

The dual variables in the transshipment problem are prices  $p(v)$  at the nodes  $v \in V$ . The dual constraints require that the *reduced cost*  $c_p(v, w) = p(v) + c(v, w) - p(w)$  is nonnegative for every arc. The complementary slackness conditions require that the reduced cost is zero for every arc with positive flow.

Orlin's algorithm uses a variant of the  $\varepsilon$ -optimality conditions that relax the flow conservation constraints rather than the complementary slackness constraints. The algorithm maintains a preflow  $f$  and a dual feasible price function  $p$  that satisfy the complementary slackness conditions. We say that such a pair is  $\Delta$ -optimal if  $|b(v) - e_f(v)| \leq \Delta$  for every node  $v$ .

First consider a simple variant of the Edmonds and Karp algorithm. One step of the algorithm selects a node  $v$  with large excess (with  $e_f(v) \gg b(v)$ ) and a node  $w$  with large deficit ( $e_f(w) \ll b(w)$ ). It sends flow from  $v$  to  $w$  along the cheapest path between them. Such a step decreases the excess or the deficit of a node. Roughly speaking the maximum excess or deficit decreases by a constant factor after every  $n$  shortest path computations. Therefore, if  $B = ||b||_\infty$ , then this algorithm finds an  $\varepsilon B$ -optimal transshipment in  $O(n \log(1/\varepsilon))$  shortest path computations. If all demands are integral, then throughout the flow is integral, and hence after at most  $O(n \log B)$  shortest path computations an optimal transshipment is found.

The constraint matrix,  $A$ , of the transshipment problem is totally unimodular, i.e.  $\Delta(A)$  is 1. Therefore, the nonnegativity requirement for an edge  $(v, w)$ , can be deleted if the value of a  $\Delta$ -optimal flow  $f$  on an arc  $(v, w)$  is at least  $n\Delta$ . Deleting the nonnegativity requirement,  $f(v, w) \geq 0$ , corresponds, in graph theoretic terms, to contracting the arc  $(v, w)$ .

The general iterative method starts each iteration by some preprocessing to guarantee that at least one inequality will be deleted each iteration. There is no need for preprocessing in the case of the transshipment problem. It is easy to see that a  $\Delta = B/n^2$ -optimal flow must have an arcs that can be contracted. Indeed, let  $v$  be the node with largest demand in absolute value. In a  $\Delta$ -optimal flow all excess or deficit is at most  $\Delta$ , hence most of the demand of  $v$  must be satisfied, and at least one arc entering or leaving  $v$  must carry sufficient flow that it can be contracted.

The above ideas give the sketch of a strongly polynomial algorithm for the transshipment problem that consists of  $O(n^2 \log n)$  shortest path computations. One iteration of the algorithm computes and  $\varepsilon$ -optimal solution for some  $\varepsilon$ . This consists of  $O(n \log n)$  shortest path computations. Then some edges are contracted (at least one each iteration), and a new iteration starts. One reason why this algorithm is slow is, that it restarts the flow computation after every iteration. Orlin [24] gave a more sophisticated version of this algorithm that does not have separate iterations, but instead contracts edges and then continues the same flow computation. Amortized analysis is used to show that a minimum cost transshipment is found in  $O(n \log n)$  shortest path computations.

Let us compare running times of the resulting strongly polynomial algorithms with the algorithm of Edmonds & Karp. Note that while the iterative algorithm

is already strongly polynomial, its running time does not compare favorably with the Edmonds & Karp algorithm. In order to make the algorithm strongly polynomial, the  $O(\log B)$  in the running time had to be replaced by  $O(m \log n)$ . However, in the running time of Orlin's algorithm the  $O(\log B)$  is replaced by only an  $O(\log n)$ , which is much more attractive.

## 4.2 The Maximum-Flow Problem

The first polynomial algorithms for the maximum-flow problem, due independently to Dinic [5] and Edmonds and Karp [6], were strongly polynomial. The preflow-push algorithm of Goldberg [10] is one of the biggest recent breakthroughs in network algorithms. Using sophisticated data-structures Goldberg and Tarjan [10] obtained an  $O(nm \log(n^2/m))$  implementation of this algorithm. This was the fastest known strongly polynomial algorithm for all values of  $n$  and  $m$  until very recently. Ahuja and Orlin [2] have developed a simple data-structure free version of the algorithm, that is more efficient for all but very large values of the capacities. Its running time is  $O(mn + n^2 \log U)$ , if we assume that the capacities are integral and at most  $U$ . This algorithm was the starting point of the strongly polynomial algorithm of Cheriyan and Hagerup [3]. Combining the ideas for converting polynomial algorithms into strongly polynomial ones, the data-structures used by Goldberg and Tarjan and a very sophisticated amortized analysis, they give a maximum-flow algorithm that runs in  $O(mn + n^2 \log^3 n)$  time. This time bound is better than the Goldberg and Tarjan bound if the graph is not too sparse.

For practical problems the size of the numbers involved often compares favorably to the size of the network. Therefore, polynomial algorithms that are not strongly polynomial might be more efficient in practice. To compare polynomial algorithms with strongly polynomial ones Gabow has suggested the *similarity assumption*, that is, to assume that all numbers involved are integers that have size at most  $O(\log n)$ . Unfortunately, algorithms whose running time depends exponentially on the size of the numbers involved appear polynomial under this assumption. A slightly weaker assumption that avoids this problem, is to assume that all numbers involved are integral of size is at most  $\log^{O(1)} n$ . Under this assumption both Orlin's transshipment algorithm and the Cheriyan & Hagerup algorithm compares favorably to its polynomial counterpart. It is an interesting open problem to see what other polynomial algorithms be converted to strongly polynomial ones by replacing the size of the numbers in the running time by  $O(\log^c n)$  for some constant  $c$  where  $n$  is the combinatorial size of the problem.

## 4.3 Multi-Commodity Flows

There are several important special cases of the linear programming problem which have a combinatorial structure similar to the transshipment problem and also play a fundamental role in applications, but are not nearly as well understood. For many of these problems, the only known polynomial algorithm is obtained by

using polynomial linear programming algorithms. Such algorithms use continuous methods and do not take advantage of the combinatorial structure. It is an interesting question whether some of these problems can be solved more efficiently using combinatorial techniques. The multi-commodity flow problem is an example of such a problem.

Many polynomial algorithms, if stated with an appropriate initial point, find an approximately optimal solution in time polynomial in  $n, m$  and  $\log(1/\varepsilon)$ . If  $\varepsilon$  is large this is strongly polynomial. In the definition of  $\varepsilon$ -optimality used for designing strongly polynomial algorithms, the complementary slackness conditions have been relaxed in a one-sided way. A symmetric version would require only  $y_i \leq \varepsilon$  instead of  $y_i = 0$ . Primal-dual potential reduction versions of Karmarkar's method find " $\varepsilon$ -optimal" solutions in this symmetric sense in time polynomial in  $n, m$  and  $\log(1/\varepsilon)$ . See Todd [31] for a survey of such methods.

Recently, Klein, Stein and Tardos [16] have obtained a combinatorial algorithm that significantly outperforms the general purpose techniques for a special case of the multi-commodity flow problem where  $\varepsilon$  is large. Consider a multi-commodity flow that satisfies the demands, but not necessarily the capacity constraints. Let  $f(v, w) = \sum_i f_i(v, w)$  denote the sum of the different commodities on an arc  $(v, w)$ . We consider the maximum ratio  $\lambda = \max_{(v,w) \in E} f(v, w)/u(v, w)$  as the value of the flow. The objective of the *concurrent flow problem* is to minimize this ratio  $\lambda$ .

The concurrent flow problem has one variable,  $\lambda$  whose coefficients are not  $0, \pm 1$  in the corresponding linear program. The coefficients of  $\lambda$  are the capacities. The problem is an extension with one new variable of a linear program that is known to be solvable in strongly polynomial time. This implies that the problem can be solved in strongly polynomial time [22].

Shahrokhi and Matula [28] gave a fully polynomial approximation algorithm for the special case of this problem with unit capacities, that is, an algorithm that finds a concurrent flow with value no worse than  $(1 + \varepsilon)$  times the optimal, in time polynomial in  $n, m$  and  $1/\varepsilon$ . This algorithm is quite slow, both in terms of its dependence on  $\varepsilon$  and in its dependence on the other parameters. Klein, Stein and Tardos [16] gave a much faster algorithm extending the Shahrokhi & Matula technique. The algorithm uses a notion of  $\varepsilon$ -optimality similar to the symmetric notion used by the general linear programming algorithms. The algorithm runs roughly (ignoring  $\log$ 's) in  $O(\varepsilon^{-2}m^2k)$  time. This compares quite favorably to known algorithms as long as  $\varepsilon$  is large (e.g., a constant).

Approximation algorithms for the concurrent flow problem are especially interesting because several of its applications require only approximately optimal solutions. Examples of such applications are graph partitioning by Leighton and Rao [17] and VLSI routing by Raghavan and Thompson [26]. It remains a significant open problem to design algorithms for this problem that outperform general purpose linear programming techniques for the case when  $\varepsilon$  is small.

The key dual variables for the concurrent flow problem are nonnegative lengths  $\ell(v, w)$  associated with the edges. Let  $\lambda$  denote the value of the multi-commodity flow  $f$ . The complementary slackness conditions stated in terms of these variables require that (1) all arcs with nonzero length are fully utilized, that

is, for each arc  $(v, w)$  either the total flow on it is exactly  $\lambda$  times its capacity, or the arc has zero length; and (2) for each commodity the flow can be decomposed into flows along simple paths each of which is shortest according to the length function  $\ell$ .

The variant of  $\epsilon$ -optimality used by the algorithm is defined by relaxing both of the above constraints. We require that (1) each arc is either almost fully utilized or has close to zero length; and (2) almost all the flow is carried on close to shortest paths. It is not too difficult to prove that these conditions imply that the value of multi-commodity flow is close to optimal.

The idea of the approximation scheme is to maintain a length function defined as a function of the flow,  $\ell(v, w) = \exp(\alpha f(v, w))$  for an appropriate constant  $\alpha$ , such that the first part of the  $\epsilon$ -optimality conditions is satisfied. The second assumption is gradually enforced as flow is repeatedly rerouted from long paths to the corresponding shortest path.

## 5. Generalized Flow: An Open Problem

The most intriguing special case of linear programming for which no strongly polynomial algorithm is known is the generalized flow problem, which is essentially the problem of an arbitrager who wants to maximize his profit by converting currencies at different exchanges. The problem is defined by a network, a source  $s \in V$ , capacities and positive gainfactors  $\gamma(v, w)$  on the arcs. When  $x$  units of flow enter an arc  $(v, w)$  then  $x\gamma(v, w)$  units leave the arc. The excess of a preflow  $f$  is defined to be  $\sum_w f(w, v)\gamma(w, v) - \sum_w f(v, w)$ . A preflow is a *generalized flow* if all nodes, except the source, have zero excess. The problem is to find a generalized flow with maximum excess at the source.

The generalized flow problem is a linear programming problem. The constraint matrix is not 0,  $\pm 1$ ; its nonzero entries are the gainfactors. The problem can also be considered as a slight generalization of the class of feasibility problems with two variables per inequality considered by Megiddo [21]. Consider the (equivalent) version of the problem with demands at the nodes instead of the capacities. The linear programming dual of this version has two variables in each inequality. However, it does not fit into Megiddo's framework, because it has an objective function with more than two variables; the demands at the nodes are the coefficients of the objective function in the dual.

The fastest known algorithm for the generalized flow problem is due to Vaidya [32]. It combines general linear programming techniques and fast matrix inversion techniques. It makes some use of the network structure by speeding up the matrix inversions involved. Goldberg, Plotkin and Tardos [8] gave two combinatorial algorithms for this problem. One of the algorithms uses a transshipment algorithm as a subroutine, the other one combines techniques from maximum-flow and transshipment algorithms. Both run in time comparable to the fastest known algorithms for the problem. It remains an interesting open problem to find algorithms that significantly outperform general linear programming techniques.

Finally, we point out that a combination of the multi-commodity flow and the generalized flow problem covers the full generality of linear programming.

We have mentioned in the first part of of Section 4. that every program can be reduced to an equivalent 2-commodity flow problem [1,12]. However, the reduction is not strongly polynomial. The following theorem is an easy extension.

**Theorem 7.** *Every linear program can be reduced in strongly polynomial time to an equivalent 2-commodity generalized flow problem.*

## References

1. G. M. Adelson-Velskii, E. A. Dinic, A. V. Karzanov: Flow algorithms. Izdatelstvo Nauka, Moskwa, 1975 (in Russian)
2. R. K. Ahuja, J. B. Orlin: A fast and simple algorithm for the maximum flow problem. *Operations Res.* **37**, 748–759 (1989)
3. J. Cheriyan, T. Hagerup: A randomized maximum-flow algorithm. In: *Proceedings of the 30st Annual Symposium on Foundation of Computer Science*, 1989.
4. W. Cook, A. M. H. Gerards, A. Schrijver, É. Tardos: Sensitivity theorems in integer and linear programming. *Math. Programming* **34**, 251–264 (1986)
5. E. A. Dinic: Algorithm for solution of a problem of maximum flow in networks with power estimation. *Sov. Math. Dokl.* **11**, 1277–1280 (1970)
6. J. Edmonds, R. M. Karp: Theoretical improvements in algorithmic efficiency for network flow problems. *J. Assoc. Comput. Mach.* **19**, 248–264 (1972)
7. A. Frank, É. Tardos: An application of simultaneous Diophantine approximation in combinatorial optimization. *Combinatorica* **7** (1) 49–65 (1987)
8. A. V. Goldberg, S. A. Plotkin, É. Tardos: Combinatorial algorithms for the generalized circulation problem. In: *Proceedings of the 29th FOCS Symposium*, 1988, pp. 432–443 (to appear in *Mathematics of Operations Research*)
9. A. V. Goldberg, R. E. Tarjan: Finding minimum-cost circulations by canceling negative cycles. *J. Assoc. Comput. Mach.* **36**, 873–886 (1989)
10. A. V. Goldberg, R. E. Tarjan: A new approach to the maximum flow problem. *J. Assoc. Comput. Mach.* **35**, 921–940 (1988)
11. M. Grötschel, L. Lovász, A. Schrijver: The ellipsoid method and its consequences in combinatorial optimization. *Combinatorica* **1**, 169–197 (1981)
12. A. Itai: Two-commodity flow. *JACM* **25**, 596–611 (1978)
13. N. Karmarkar: A new polynomial-time algorithm for linear programming. *Combinatorica* **4**, 373–395 (1984)
14. R. M. Karp, C. H. Papadimitriou: On linear characterization of combinatorial optimization problems. *SIAM J. Comput.* **11**, 620–632 (1982)
15. L. G. Khachyian: A polynomial algorithm in linear programming. *Dokl. Akad. Nauk SSSR* **244**, 1093–1096 (1979) [English translation: *Sov. Math. Dokl.* **20** (1979) 191–194]
16. P. Klein, C. Stein, É. Tardos: Leighton-Rao might be practical: Faster approximation algorithms for concurrent flow with uniform capacities. In: *Proceedings of the ACM Symposium on the Theory of Computing*, 1990, pp. 310–321
17. T. Leighton, S. Rao: An approximate max-flow min-cut theorem for uniform multi-commodity flow problems with applications to approximation algorithms. In: *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, 1988, pp. 422–431
18. A. K. Lenstra, Jr. H. W. Lenstra, L. Lovász: Factoring polynomials with rational coefficients. *Math. Ann.* **261**, 515–548 (1982)
19. O. L. Mangasarian: Condition number for linear inequalities and linear programs. *Methods of Operations Research* **43**, 3–15 (1981)

20. N. Megiddo: Combinatorial optimization with rational objective functions. *Mathematics of Operations Research* **4**, 414–424 (1979)
21. N. Megiddo: Towards a genuinely polynomial algorithm for linear programming. *SIAM J. Comput.* **12**, 347–353 (1983)
22. C. H. Norton, S. Plotkin, É. Tardos: Using separation algorithms in fixed dimension. In: *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms 1990*, pp. 377–387
23. J. B. Orlin: On the simplex algorithm for networks and generalized networks. *Math. Prog. Studies* **24**, 166–178 (1985)
24. J. B. Orlin: A faster strongly polynomial minimum cost flow algorithm. In: *Proc. 20th ACM Symp. on Theory of Computing 1988*, pp. 377–387
25. M. W. Padberg, M. R. Rao: The Russian method for linear programming III: Bounded integer programming. Technical Report 81-38, New York University, Graduate School of Business Administration, 1981
26. P. Raghavan, C. D. Thompson: Randomized rounding: a technique for provably good algorithms and algorithmic proofs. *Combinatorica* **7**, 365–374 (1987)
27. A. Schrijver: *Theory of integer and linear programming*. Wiley and Sons, 1986
28. F. Shahrokhi, D. W. Matula: The maximum concurrent flow problem. *J. ACM* **37** (2), 318–334 (1990)
29. É. Tardos: A strongly polynomial minimum cost circulation algorithm. *Combinatorica* **5**(3) 247–255 (1985)
30. É. Tardos: A strongly polynomial algorithm for solving combinatorial linear programs. *Operations Research*, pp. 250–256 (1986)
31. M. J. Todd: Recent developments and new directions in linear programming. In: M. Iri, K. Tanabe (eds.) *Mathematical Programming: Recent Developments and Applications*. KKT Scientific Publishers, Tokyo 1989, pp. 109–158
32. P. M. Vaidya: Speeding up linear programming using fast matrix multiplication. In: *Proc. 30th IEEE Annual Symposium on Foundations of Computer Science*, 1989
33. D. B. Yudin, A. S. Nemirovskii: Evaluation of the informational complexity of mathematical programming problems. *Ekonomika i Matematicheskie Metody* **12**, 128–142 (1976) (in Russian)

# Computational Complexity in Finite Groups

László Babai \*

Department of Computer Science, University of Chicago, Chicago, IL 60637, USA  
and Department of Algebra, Eötvös University, Budapest, Hungary H-1088

We survey recent results on the *asymptotic* complexity of some of the fundamental computational tasks in finite groups in a variety of computational models. A striking recent feature is that techniques motivated by the problems of the more abstract models (nondeterminism, extreme parallelization) have turned out to provide powerful tools in the design of surprisingly efficient algorithms on realistic models (e.g. a *nearly linear time* membership test for permutation groups with a small base).

The techniques involve a combination of *elementary combinatorial results* on finite groups, some classical elementary group theory, and the extensive use of certain consequences of the *classification of finite simple groups* (CFSG).

Most of the recent work surveyed is due to E. M. Luks, G. Cooperman, L. Finkelstein, Á. Seress, E. Szemerédi, and the author.

## 1. Group Models and Measures of Complexity

Rubik's Cube illustrates some of the basic problems of computational group theory. We may want to decide whether or not a particular configuration is feasible (accessible without pulling the cube apart); determine the total number of feasible configurations; or construct "typical" configurations. In group theoretical terms, we are given a group  $G$  by a list  $S$  of generators (the "legal moves"), and we wish to determine whether or not a particular element of a larger group belongs to  $G$  (*membership testing*); determine the *order* of  $G$ ; generate uniformly distributed *random members* of  $G$ . The gourmet will ask more sophisticated questions such as deciding solvability, nilpotence, constructing normal closures, the center, composition factors, Sylow subgroups, etc.

The *cost* of answering these questions depends on two factors: the way group operations are performed, and the measure of cost.

For greatest generality, we consider *black box groups*, a model where no restriction is made on the way group operations are performed. In this model, elements of an unknown group  $B$  (the "group in the box") are encoded by

---

\* Research supported in part by NSF Grant CCR-8710078 and Hungarian National Foundation for Scientific Research Grant 1812.

binary strings of uniform length  $n$ . (In particular,  $|B| \leq 2^n$ .) Group operations are performed by the “black box” at unit cost. A *black box group* is a subgroup of  $B$  given by a list of generators. (Generators of  $B$ , or a recognition method of strings in  $B$ , are *not* assumed to be known.)

Two implementations of the “black box” are of particular interest: *permutation groups* (subgroups of the symmetric group  $B = \text{Sym}(\Omega)$ ), and *matrix groups* over finite fields (subgroups of the linear group  $B = GL(d, q)$ ).

The *models of computation* to be studied include deterministic as well as Monte Carlo (randomized) computation, both sequential and parallel (several processors). Timing estimates refer to the logarithmic cost RAM model [AHU].

A *Monte Carlo algorithm* uses randomization, hence its outcome may be in error. However, on any input, the probability of error is required to be  $\leq 1/4$ . By repeating the algorithm  $m$  times and taking majority vote, the chance of error is reduced to  $< e^{-m/8}$  (Chernoff’s bound).

The *cost* (in terms of a specific resource such as time, space, number of processors, length of proofs) is measured as a function of the *length of the input*, i.e. the number of input bits. An algorithm is said to have cost  $O(f(n))$  (or cost  $O^\sim(f(n))$ ) if for  $n \geq n_0$  and on *all* inputs of length  $n$ , the cost is at most  $cf(n)$  (at most  $f(n)(\log n)^c$ , resp.) ( $n_0, c$  will denote various constants throughout).

A function  $f(n)$  is *polynomially bounded* (or “*short*”, “*small*”) if  $f(n) \leq n^c$  for some  $c$  and all  $n \geq n_0$ . A *polylog bound* means  $f(n) \leq (\log n)^c$ . When used as technical terms, “*short*”, “*small*” will be Italicized.

*NC* (“Nick’s Class”) denotes (somewhat informally) the class of functions computable in *polylog* time, using a *small* number of parallel processors.

*NP* (“nondeterministic polynomial time”) stands for the class of decision problems where the “yes” answers have *short* proofs. More precisely, a predicate  $A(x)$  belongs to *NP* if there exists a polynomial time computable predicate  $B(x, w)$  such that for every input string  $x$ ,  $A(x) \leftrightarrow (\exists^P w)B(x, w)$ , where  $\exists^P w$  refers to *short* strings  $w$ . The string  $w$  is called a *witness* of the statement  $A(x)$ . The negations of *NP*-predicates form the class *coNP*. (Cf. [GJ].)

We shall also consider the class *AM* (“Arthur–Merlin”), a randomized extension of *NP*, defined as follows: the predicate  $A(x)$  belongs to *AM* if there exists an *NP*-predicate  $B(x, r)$  such that for every input string  $x$ ,  $A(x)$  is equivalent to  $B(x, r)$  for most *short* strings  $r$ . (The definition of “most” is flexible; asking more than 51% will define the same complexity class as asking, say, a  $1 - 2^{-n}$  fraction, where  $n$  is the length of  $x$ .) Informally,  $A(x)$  has *short* “interactive proofs” in the sense that if the all-knowing but untrusted Merlin is able to present a *short* “witness” in response to a random question  $r$  of polynomial time bounded Arthur, this should convince skeptical Arthur by way of overwhelming statistical evidence that Merlin’s claim  $A(x)$  is true. (Cf. [Ba2, BM, GMR, Go].)

## 2. General Methods: Black Box Groups

In this section we demonstrate the somewhat unexpected fact that nontrivial computational tasks, such as constructing *random elements* and deciding *solvability*, can be accomplished in Monte Carlo polynomial time in the extremely

general model of *black box groups*. Here and throughout the paper,  $G$  will denote a finite group and  $S$  a set of generators of  $G$ .

## 2.1 Two Combinatorial Lemmas

We begin with two elementary results which will play a key role. The first one concerns the number of group operations required in order to construct an element from a given set of generators.

A *straight line program* from  $S \subseteq G$  to  $g \in G$  is a sequence  $u_1, \dots, u_m$  of elements of the group  $G$  such that each  $u_i$  either belongs to  $S$  or is obtained from one or two previous elements by a group operation; and  $g = u_m$ . The *straight line cost*  $\text{cost}(g|S)$  is the smallest  $m$  such that such a straight line program exists.

**Lemma 1 (Reachability Lemma [BSz]).** *Given any set  $S$  of generators of a group  $G$  and any  $g \in G$ , we have  $\text{cost}(g|S) \leq (1 + \log |G|)^2$ .*

(All  $\log$ 's in this paper are to the base 2.) A *subproduct* of the elements  $h_1, \dots, h_k \in G$  is a product of the form  $h_1^{e_1} \cdots h_k^{e_k}$ , where  $e_i \in \{0, 1\}$ . The  $k$ -dimensional *cube*  $C = C(h_1, \dots, h_k)$  is the set of all the  $2^k$  subproducts of the  $h_i$ . The cube  $C$  is *nondegenerate* if  $|C| = 2^k$ . The basic structure established in the proof of Lemma 1 is a *chain of nondegenerate cubes*: we prove the existence of a sequence of elements  $h_1, \dots, h_t$  which generate a nondegenerate cube  $C$  such that  $G = C^{-1}C$ ; and for every  $i$ ,  $\text{cost}(h_i|h_1, \dots, h_{i-1}, S) \leq 2i - 1$ . The  $h_i$  are found inductively; we can continue as long as  $C^{-1}C \neq G$ : the element outside  $C^{-1}C$  of lowest straight line cost will do. We shall refer to this procedure as *doubling the cube*. Clearly, we must stop at some  $t \leq \log |G|$ .

The proof just sketched is non-constructive; it does not tell how to find an element that will double the current cube. The following lemma provides the key to an efficient Monte Carlo procedure.

A *graph* is a pair  $X = (V, E)$  where  $V$  is the set of *vertices*, and  $E$  is a set of unordered pairs of vertices, called *edges*. Two vertices  $v, v'$  are *adjacent* if  $\{v, v'\} \in E$ . An *isomorphism* of two graphs is a bijection of the vertex sets preserving adjacency. The group of self-isomorphisms of  $X$  is the *automorphism group*  $\text{Aut}(X)$ . We say that  $X$  is *vertex-transitive* if  $\text{Aut}(X)$  is a transitive subgroup of  $\text{Sym}(V)$ . The number of vertices adjacent to  $v \in V$  is the *degree* of  $v$ .  $X$  is *locally finite* if its vertices have finite degrees. The boundary  $\partial W$  of a subset  $W \subseteq V$  consists of all vertices in  $V \setminus W$  adjacent to some vertex in  $W$ .  $X$  is *connected* if  $\partial W \neq \emptyset$  for any nonempty proper subset  $W$ . A walk of length  $\ell$  in  $X$  is a chain of  $\ell + 1$  vertices, each adjacent to its predecessor. The *distance* of  $v, v' \in V$  is the length of the shortest walk between them. Let  $X'(v)$  denote the set of vertices at distance  $\leq t$  from  $v$ .

**Lemma 2 (Local Expansion Lemma [Ba4,5]).** *Let  $X = (V, E)$  be a locally finite connected vertex-transitive graph and  $v \in V$ . If  $W \subseteq X'(v)$  and  $|W| \leq |V|/2$  then  $|\partial W| \geq |W|/(4t)$ .*

This lemma has the interesting consequence that random walks on a vertex-transitive graph “don't get stuck” in a corner.

**Theorem 3** [Ba5]. *Let  $X$  be a connected vertex-transitive graph of finite degree  $d$ . Assume that  $|X^{4t}(v)| \leq |V|/2$ . Let  $\tau$  be a random integer chosen uniformly from  $\{t, t+1, \dots, \ell\}$  where  $\ell = Ct^2d \log |G|$  ( $C = 500$ ). Then with probability  $\geq 1/16$ , a random walk of length  $\tau$ , starting at  $v$ , will end outside  $X^t(v)$ .*

Among the ingredients of the proof is a Cheeger-type [Ch] eigenvalue estimate for graphs, derived from the local expansion property, following the lines of [Alo].

The graphs this theorem will be applied to are *Cayley graphs*. The vertex set of the Cayley graph  $X(G, S)$  is  $G$ ; and  $g \in G$  is adjacent to  $gh$  for  $h \in S \cup S^{-1}$ .

## 2.2 Membership and Random Generation

We begin with a nondeterministic result.

**Theorem 4** [BSz]. *Membership in black box groups belongs to NP.*

Indeed, a *short* straight line program qualifies as a *witness* of membership.

There is little hope for making this proof constructive, even in the very special case when the “group in the box” is the multiplicative group of  $GF(q)$ . Indeed in this case, finding a straight line program to generate  $g$  from  $S = \{h\}$  is equivalent to solving the equation  $h^x = g$ . This is the *discrete logarithm* problem, not believed to be solvable in polynomial time. (Known algorithms require time  $\exp(c\sqrt{q \log q})$ ; whereas polynomial time would mean  $\text{polylog}(q)$  steps.)

Yet, part of the proof *can* be turned into an efficient Monte Carlo algorithm.

**Theorem 5** [Ba5]. *Nearly uniformly distributed random elements of a black box group can be constructed in Monte Carlo polynomial time.*

Nearly uniform distribution means each element has probability  $(1 \pm \epsilon)/|G|$  to be selected; and the reliability of the algorithm is  $\geq 1 - \delta$ , where  $\epsilon, \delta$  are input parameters, and the number of operations is polynomially bounded in  $k = \log |G| + \log(1/\epsilon) + \log(1/\delta)$ . The idea is to construct a set  $S'$  of generators such that the diameter of the Cayley graph  $X(G, S')$  is *small*. Once this is accomplished, *short* random walks are known to produce nearly uniformly distributed elements [Ald, Alo].

To reduce the diameter, we should like to adapt the “doubling the cube” trick. The difficulty is, how to obtain the next  $h_i$ , which must be outside the set  $C^{-1}C$ , where  $C$  is the current cube. The solution is a *short* random walk of random length over the Cayley graph. By Theorem 3, such a walk has a fixed positive chance of reaching a desired element.

## 2.3 Nonmembership, Order

These two problems are known to belong to the class *AM* [Ba4]. Indeed, to prove this was the original motivation behind inventing the Local Expansion Lemma.

For black box groups, the nonmembership and order verification problems are provably *not* in *NP* [BSz]. (To be precise, here we are talking about a *relativized* version of *NP*: computations refer to the black box, an “oracle” [GJ].)

We conjecture, however, that for matrix groups over finite fields, these problems belong to  $NP$ . This will follow from the conjecture below. The *length* of a *presentation* of a group (in terms of generators and relations) is the number of bits required to write down the presentation. E. g., the presentation  $\langle a | a^N = 1 \rangle$  of the cyclic group has length  $\log N + O(1)$  (we write exponents in binary).

**Short Presentation Conjecture.** (i) Every finite simple group  $G$  has a presentation  $R$  of length  $\text{polylog}(|G|)$ . (ii) If  $G$  is of Lie type over  $GF(q)$  then such  $R$  is computable from the standard name of  $G$  in time  $\text{polylog}(|G|)$ , assuming  $GF(q)$  and a primitive root in it are explicitly given.

One can prove, using Lemma 1, that part (i) of the Conjecture, if true, automatically extends to all finite groups [BKLP]. The conjecture itself has been verified for all  $G$  except those of rank one twisted Lie type [BKLP], cf. [Ka4]. Note that for a Lie-type simple group of rank  $d$  over  $GF(q)$ , the Conjecture requires presentations of length  $\leq (d \log q)^c$ ; whereas the Steinberg presentations [St, Car] require an exponentially greater number, about  $d^2 q$  generators.

Verification of the *order* is a central problem. If it belongs to  $NP$  (as we expect for matrix groups), this brings a number of other verification problems into  $NP$ , including composition factors, homomorphisms, isomorphisms, kernels, minimal normal subgroups. On the other hand, problems known to be in  $AM$  (for black box groups in general) but not expected to belong to  $NP$  (even for permutation groups) include verification of the intersection of two subgroups, the centralizer of an element, non-conjugacy of two elements [Ba4].

## 2.4 Random Subproducts

Algorithms often depend on access to random elements of  $G$  (e.g. [CFS, NP]). Theorem 5 constructs such elements in reasonable polynomial time, but not efficiently enough for some applications. *Random subproducts*, however, often emulate truly random elements very efficiently.

Let  $S = \{g_1, \dots, g_s\}$ . A *random subproduct* is a subproduct  $\varphi = g_1^{e_1} \cdots g_s^{e_s}$ , where the  $e_i$  are independent uniform  $(0, 1)$ -variables (coin flips).

**Lemma 6** [BLS2]. *Let  $H < G$  be a proper subgroup. Then  $\text{Prob}(\varphi \notin H) \geq 1/2$ .*

Let  $L$  be the maximum length of subgroup chains in  $G$ . ( $L \leq \log |G|$ ) Lemma 6 implies that the probability that  $2L(1+\alpha)$  random subproducts do not generate  $G$  is less than  $\exp(-\alpha^2 L / (1 + \alpha))$ . A refined argument yields:

**Lemma 7** [BCFLS]. *If  $G$  is given by a list of  $s$  generators, then a Monte Carlo procedure, using  $O(s \log L)$  group operations, produces (with large but fixed probability) a set of  $O(L)$  generators for  $G$ .*

This keeps the number of generators down when constructing subgroups. Some additional combinatorics yields a particularly efficient *normal closure* algorithm. Recall that the notation  $O^\sim(f(n))$  refers to an upper bound  $f(n)(\log n)^{O(1)}$ .

**Theorem 8** [BCFLS]. *Let  $H \leq G$  be black box groups, each given by a list of  $O(L)$  generators. Then the normal closure of  $H$  in  $G$  can be constructed (in the form of  $O(L)$  generators) in Monte Carlo  $O^{\sim}(L^2)$  operations.*

As an immediate application, we obtain polynomial time Monte Carlo algorithms to decide *solvability* and *nilpotence* of  $G$  (both require  $O^{\sim}(L^3)$  operations).

We should stress that the results of this section apply to all black box groups, including their implementations as *matrix groups* or *permutation groups*. In spite of their generality, the results are strong enough to yield asymptotic savings even in the well-studied area of permutation groups (see Sect. 4.2).

### 3. Permutation Groups: Survey of Complexity Status

We consider groups  $G \leq \text{Sym}(\Omega)$  where  $|\Omega| = n$ . The *stabilizer* of  $x \in \Omega$  is the subgroup  $G_x = \{g \in G : x^g = x\}$ . The *pointwise set stabilizer* of  $\Delta \subseteq \Omega$  is the subgroup  $G_\Delta = \bigcap_{x \in \Delta} G_x$ . We call  $\Delta$  a *base* if  $G_\Delta = \{1\}$ .

Let  $G = G^{(0)} \geq G^{(1)} \geq \dots \geq G^{(m)} = \{1\}$  be a chain of subgroups. A *strong generating set* (SGS) w.r. to this chain is a set  $S \subseteq G$  such that  $G^{(i)} = \langle S \cap G^{(i)} \rangle$  for every  $i$ . A *transversal system* (TS) is a family  $\{T_i : 1 \leq i \leq m\}$ , where  $T_i$  is a (right) transversal (set of coset representatives) of  $G^{(i)}$  in  $G^{(i-1)}$ . A *partial transversal system* is a family of partial transversals  $T'_i \subseteq T_i$ .

The *stabilizer chain* w.r. to a given ordered base  $\Delta = \{x_1, \dots, x_m\}$  is defined as  $G^{(i)} = G_{x_1, \dots, x_i}$ . The concept of an SGS w.r. to an ordered base was introduced by C. C. Sims in the early 60's as a central tool in computational group theory ([Sim1, 2]). Given an SGS, a TS is readily constructed, solving the *membership* and *order* problems and the construction of truly uniformly distributed *random elements*. A slight modification yields *normal closures*, clearing the way for more advanced applications. The central problem of constructing an SGS was efficiently solved by Sims [Sim1, 2].

The *asymptotic complexity* of these algorithms was not analyzed until 1980 when the complexity of many of these algorithms was recognized to be polynomial time in [FHL]. In particular, an  $O(n^6 + sn^2)$  variant of Sims's SGS algorithm was constructed, where  $s$  is the number of input generators. As a consequence, membership, order, normal closures, solvability were shown to be computable in polynomial time [FHL]. E. M. Luks has subsequently added an array of elegant polynomial time algorithms which, for the first time, required deeper group theoretic analysis. The list includes the center, a composition chain [Lu2], and subcases of the *coset intersection* problem, i.e. determining  $G \cap Hh$  where  $G, H \leq \text{Sym}(\Omega)$ ,  $h \in \text{Sym}(\Omega)$ . The subcases solved in polynomial time in Luks's seminal paper [Lu1] include the case when  $G$  is solvable, or more generally, the nonabelian composition factors of  $G$  are restricted to the set  $\mathcal{G}(c)$  consisting of the alternating groups of degree  $\leq c$ , the groups of Lie type of rank  $\leq c$ , and the sporadic groups. The algorithm uses classical *divide and conquer* algorithmic techniques, splitting the domain into orbits and then into domains of imprimitivity ([Wi]). When  $G$  is primitive, some of the algorithms use exhaustive search. In such cases, the polynomial time claim depends on the following result.

**Theorem 9** [BCP]. *If  $G \leq \text{Sym}(\Omega)$  is primitive and  $G \in \mathcal{G}(c)$  then  $G$  is small.*

(Small means  $|G| \leq n^{c'}$  for some constant  $c'$ , depending on  $c$ .) For primitive solvable groups, the precise bound is  $|G| \leq 24^{-1/3}n^c$  where  $c = 1 + \log_9(48 \cdot 24^{1/3}) = 3.24399\dots$  [Pá, Wo].

Sylow subgroups and Sylow normalizers were added to the polynomial time library by Kantor [Ka1, 2, 3]. For a long list of additional results see [KL].

Another important observation of [Lu1] was that a number of problems, including coset intersection, setwise stabilizer of a subset, centralizer of an element, and centralizer of a subgroup, are *equivalent* (polynomial time reducible to one another), and the *graph isomorphism* problem (to decide whether or not two given graphs are isomorphic) is reducible to each. In particular, as long as graph isomorphism is not solved in polynomial time (the best current algorithm requires  $\exp(O^{\sim}(\sqrt{n}))$  for graphs on  $n$  vertices, cf. [BL]), coset intersection, etc., are not expected to be efficiently solvable. On the other hand, the decision versions of these problems (“is  $G \cap Hh \neq \emptyset$ ?”) are *not NP-complete*, unless the so-called *polynomial time hierarchy* of complexity classes collapses [BM, GMW]. (The conjecture that the polynomial time hierarchy does not collapse is a stronger version of the famous  $NP \neq coNP$  conjecture [Sto].)

Some related problems are *NP-complete*; the nicest is A. Lubiw’s result: the predicate “ $G$  has a fixed-point-free element” is *NP-complete*, even for elementary abelian 2-groups [Lub]. An even harder problem is to determine minimum generating sequences; the length of the shortest word in  $S$  representing  $g \in G$  is *PSPACE-complete* [Je1]. (For related problems, see [BHKLS].)

On the other end of the spectrum, some of the basic problems were shown to admit ultra-fast *parallel algorithms*. Most notably, *membership*, *order*, and even a *composition chain* are computable in *NC* [BLS1]. A striking feature of the algorithm is that even for the rudimentary tasks of membership testing, we are forced to determine the composition factors first, using several facts of asymptotic group theory currently derivable only via the *classification of finite simple groups* (CFSG) (cf. Sect. 5). – Coset intersection is *not* known to be in *NC*; if it is in *NC*, then so is the isomorphism of graphs of degree 3 [LM].

## 4. Efficient Construction of Strong Generators

The  $O(n^6 + sn^2)$  analysis of the SGS algorithm of [FHL] was soon replaced by  $O(n^5 + sn^2)$  [Kn, Je2]. Knuth’s is the closest to Sims’s original approach and is quite efficient in practice, but there exist large collections of examples where its typical behavior is as bad as the  $n^5$  worst case bound (while  $s = n - 1$ ) [Kn].

The  $n^5$  bottleneck was broken, using machinery developed for the *NC* result, in [BLS2] ( $O^{\sim}(n^4 + sn^2)$ ). Further improvements yield  $O^{\sim}(sn^3)$  [BLS3], the best deterministic bound to date. These results heavily depend on the CFSG.

### 4.1 A Fast and Elementary Monte Carlo SGS Algorithm

We outline a new  $O^{\sim}(n^3 + sn)$  time Monte Carlo SGS algorithm with a perfectly *elementary* analysis [BCFLS].

Let  $H \leq G$  be a subgroup of index  $k$ ; and  $T = \{t_1, \dots, t_k\}$  a right transversal. For  $g \in G$ , let  $\bar{g} = t_i$  where  $Hg = Ht_i$ . Schreier's lemma asserts that the  $sk$  elements  $t_i h_j \overline{t_i h_j}^{-1}$  ("Schreier generators") generate  $H$ , where  $S = \{h_1, \dots, h_s\}$  [Ha]. Noting that a transversal for  $G_x$  in  $G$  is easily constructed, a natural approach to constructing an SGS would be to consider the Schreier generators for  $G_x$  and repeat. The difficulty is that the number of generators grows rapidly.

*Random subproducts* are the new tool. (When using the results of Sect. 2.4, the following bound comes handy: *Every subgroup chain in  $S_n$  has length  $< 2n$*  [Ba3, CST].) Lemma 7 alone saves nearly an order of magnitude over [Kn, Je2]: we reduce the number of generators of  $G$  to  $O(n)$ , then construct a transversal and the Schreier generators of  $G_x$ ; repeat. The cost of this naive approach is  $O^\sim(n^4 + sn)$ . When  $G \geq A_n$ , we have a particularly speedy variant, based on the following consequence of Lemma 6.

**Proposition 10** [BLS2]. *If  $G = \langle S \rangle \leq S_n$  and  $S'$  is a set of  $c \log n$  random subproducts of  $S$  then with large probability, the orbits of  $\langle S' \rangle$  and  $G$  agree.*

Applying this to the action of  $G \geq A_n$  on the set of ordered 6-tuples we find that  $O(\log n)$  random subproducts are likely to generate all of  $G$ . So the previous argument, using only  $O(\log n)$  random subproducts in each round, constructs an SGS for  $G \geq A_n$  in time  $O^\sim(n^3)$ . This process works also when  $G$  induces  $A_k$  or  $S_k$  on some orbit  $\Delta \subseteq \Omega$  ( $|\Delta| = k$ ). However, now we don't get generators of  $G_\Delta$ . (The procedure preserves the action of  $G$  on  $\Delta$  only.) Instead, we use a set of  $O(k)$  defining relations of  $A_k$  or  $S_k$  to construct *normal generators* of  $G_\Delta$ ; and then use our *normal closure* algorithm (Theorem 8) to obtain  $G_\Delta$ . Another ingredient of the  $O^\sim(n^3 + sn)$  algorithm is *computation of  $G_x$  in  $O^\sim(sn)$* , where  $k$  is the length of the orbit of  $x$  (apply Lemma 7 to the Schreier generators). This bound is exploited through a "smallest orbit first" strategy. By adding action on maximal blocks to  $\Omega$  we ensure that the next  $x$  to be stabilized is from an orbit with primitive action. The timing depends on a combinatorial observation:

**Lemma 11** [BCFLS]. *If  $G \leq S_n$  is primitive and  $G_x$  has a nontrivial orbit of length  $k < n/2$  then every subgroup  $H \neq \{1\}$  of  $G_x$  has a nontrivial orbit of length  $\leq k$ .*

## 4.2 Small Base Groups in Nearly Linear Time

Groups with a small base are of particular importance; e.g. linear groups, treated as permutation groups on a vector space, always have a base of size  $\leq \log n$ . Let us say that a family of groups has small bases if they have base size  $\text{polylog}(n)$ . The SGS methods of [Si1,2], [Kn], [Je2], [BCFLS] require  $> n^2$  time for such groups. Combinatorial techniques based on Lemmas 1 and 2 have recently led to a *Monte Carlo SGS algorithm in time  $O^\sim(n)$  for small base groups* [BCFS].

The basic ideas are (i) a very efficient implementation of Sims's "Schreier vector" data structure to store transversals, based on the "doubling the cube" trick (Sect. 2.1); and (ii) the use of Lemma 2 to rapidly locate elements not yet reached by the current partial transversal system. A key new feature of these methods

is that rather than operating with the coarse *subgroup* structure, we are able to handle *chains of certain subsets*, such as cubes and their generalizations.

## 5. CFSG vs. Elementary

We mention some of the consequences of the simple groups classification (CFSG) used in the analysis of the algorithms quoted. *Schreier's conjecture* that the outer automorphism groups of simple groups are solvable, is used in Luks's composition chain algorithm [Lu2] and the algorithms building on it [BLS1, 2, 3, Ka1, 2, 3]. In Sect. 4.1 we used that the degree of transitivity of  $G \leq S_n$  is  $t \leq 6$  (unless  $G \geq A_n$ ) [CKS]. At the cost of some extra log factors (swallowed by the  $O^\sim$  notation) this can be replaced by the 19th century bound  $t = o(\log^2 n)$  [Jo]. Using the CFSG, Cameron has shown that if  $G$  is a primitive group of order  $> n^{2\log n}$  then  $n = \binom{k}{\ell}^m$  and  $G$  is a subgroup of  $S_k \wr S_m$  with socle  $A_k^m$  acting on the ordered  $m$ -tuples of  $\ell$ -subsets of a  $k$ -set [Cam]. This result helps reduce the case of "large" primitive groups to  $S_n$ ; the remaining primitive groups have small bases. This is indispensable for [BLS1, 2, 3], even if all we need is to *test membership!* Kantor uses detailed knowledge of the CFSG even just to find an element of order  $p$ .

Other elementary estimates that may help avoid CFSG references (at a cost of some extra log's) include the bound  $|G| < \exp(4\sqrt{n} \log^2 n)$  for  $G$  primitive but not doubly transitive [Ba1] and  $|G| < n^{c\log^2 n}$  for  $G \not\cong A_n$  doubly transitive [Py]. Bochert's 1892 estimate [Bo] that a doubly transitive group  $G \not\cong A_n$  has *minimal degree*  $\geq n/4$  (cf. [Wi]) is used in [BCFS]. *Combinatorial proofs* may directly suggest *efficient algorithms*. A case in point is the algorithm derived from a simple proof of Jordan's  $o(\log^2 n)$  bound on the degree of transitivity [BS], allowing ultra-parallelized (NC) management of  $S_n$  [BLS1].

**Conclusion.** During the past decade, the asymptotic complexity of computation in finite groups has been analyzed in a variety of models of computation. New combinatorial and algebraic tools have been developed. Structural insights gained from the study of models ranging from the *unrealistic* (extreme parallelization: NC) to the *absurd* (nondeterminism: NP, AM) have contributed to the design of new efficient algorithms with a reasonable expectation of competitive implementations.

*Acknowledgment.* I feel privileged to have had an exciting ongoing collaboration with Gene Luks, now for over a decade. I have also greatly benefited from joint work with Gene Cooperman, Larry Finkelstein, and Ákos Seress.

## References

- [AHU] A. Aho, J. Hopcroft, J. Ullman: The design and analysis of computer algorithms. Addison-Wesley 1974
- [Ald] D. Aldous: On the Markov chain simulation method for uniform combinatorial distributions etc. *Probab. Eng. Info. Sci.* **1** (1987) 33–46
- [Al] N. Alon: Eigenvalues and expanders. *Combinatorica* **6** (1986) 83–96

- [Ba1] L. Babai: On the order of uniprimitive permutation groups. *Ann. Math.* **113** (1981) 553–568
- [Ba2] L. Babai: Trading group theory for randomness. *Proc. 17th ACM STOC*, Providence RI, 1985, pp. 421–429
- [Ba3] L. Babai: On the length of chains of subgroups in the symmetric group. *Comm. Algebra* **14** (1986) 1729–1736
- [Ba4] L. Babai: Bounded-round interactive proofs in finite groups. *SIAM J. Discr. Math.* (to appear)
- [Ba5] L. Babai: Local expansion of vertex-transitive graphs and random generation in finite groups. *Proc. 23rd ACM STOC* 1991, pp. 164–174
- [BCFLS] L. Babai, G. Cooperman, L. Finkelstein, E.M. Luks, Á. Seress: Fast Monte Carlo algorithms for permutation groups. *23rd ACM STOC* 1991, pp. 90–100
- [BCFS] L. Babai, G. Cooperman, L. Finkelstein, Á. Seress: Permutation groups with small base in almost linear time. In: *Proc. ISSAC '91*, Bonn (to appear)
- [BCP] L. Babai, P.J. Cameron, P.P. Pálfy: On the orders of primitive groups with restricted nonabelian composition factors. *J. Algebra* **79** (1982) 161–168
- [BHKLS] L. Babai, G. Hetyei, W. M. Kantor, A. Lubotzky, Á. Seress: On the diameter of finite groups. *Proc. 31st IEEE FOCS* 1990, pp. 857–865
- [BKLP] L. Babai, W. M. Kantor, E. M. Luks, P. P. Pálfy: Short presentations for finite groups. In preparation
- [BL] L. Babai, E.M. Luks: Canonical labeling of graphs. *Proc. 15th ACM STOC* 1983, pp. 171–183
- [BLS1] L. Babai, E. Luks, Á. Seress: Permutation Groups in NC. *Proc. 19th ACM STOC* 1987, pp. 409–420
- [BLS2] L. Babai, E. Luks, Á. Seress: Fast management of permutation groups. *Proc. 28th IEEE FOCS* 1988, pp. 272–282
- [BLS3] L. Babai, E. Luks, Á. Seress: Fast deterministic management of permutation groups. In preparation
- [BM] L. Babai, S. Moran: Arthur–Merlin games: a randomized proof system, and a hierarchy of complexity classes. *J. Comp. Sys. Sci.* **36** (1988) 254–276
- [Bo] A. Bochert: Ueber die Classe der Transitiven Substitutionsgruppen. *Math. Ann.* **40** (1892) 176–193
- [BS] L. Babai, Á. Seress: On the degree of transitivity of permutation groups: A short proof. *J. Comb. Theory A* **45** (1987) 310–315
- [BSz] L. Babai, E. Szemerédi: On the complexity of matrix group problems I. *Proc. 25th IEEE FOCS*, Palm Beach, FL, 1984, pp. 229–240
- [Cam] P.J. Cameron: Finite permutation groups and finite simple groups. *Bull. London Math Soc.* **13** (1981) 1–22
- [Car] R. Carter: Simple groups of Lie type. Wiley 1972, 1989
- [Ch] J. Cheeger: A lower bound for the smallest eigenvalue of the Laplacian. In: *Problems in Analysis*. Princeton Univ. Press 1970, pp. 195–199
- [CFS] G. Cooperman, L. Finkelstein, N. Sarawagi: A random base change algorithm for permutation groups. In: *Proc. ISSAC 1990*, Tokyo, ACM Press and Addison-Wesley 1990, pp. 161–168
- [CKS] C.W. Curtis, W.M. Kantor, G. Seitz: The 2-transitive permutation representations of the finite Chevalley groups. *Trans. A.M.S.* **218** (1976) 1–57
- [CST] P.J. Cameron, R. Solomon, A. Turull: Chains of subgroups in symmetric groups. *J. Algebra* **127** (1989) 340–352
- [FHL] M. L. Furst, J. Hopcroft, E. M. Luks: Polynomial-time algorithms for permutation groups. *21st IEEE FOCS* 1980, pp. 36–41
- [GJ] M. Garey, D.S. Johnson: Computers and Intractability: A guide to the theory of NP-completeness. Freeman, New York 1979

- [GMR] S. Goldwasser, S. Micali, C. Rackoff: The knowledge complexity of interactive proofs. *SIAM J. Comp.* **18** (1989) 186–208
- [GMW] O. Goldreich, S. Micali, A. Wigderson: Proofs that yield nothing but their validity and a methodology of cryptographic protocol design. *Proc. 27th IEEE FOCS* 1986, pp. 174–187
- [Go] S. Goldwasser: Interactive proofs and applications. These proceedings, pp. 1521–1535
- [Ha] M. Hall, Jr.: *The theory of groups*. Macmillan, New York 1959
- [Je1] M. R. Jerrum: The complexity of finding minimum length generator sequences. *Theor. Comp. Sci.* **36** (1985) 265–289
- [Je2] M.R. Jerrum: A compact representation for permutation groups. *J. Algorithms* **7** (1986) 60–78
- [Jo] C. Jordan: Nouvelles recherches sur la limite de transitivité des groupes non alternés. *Bull. Soc. Math. France* **1** (1873) 35–60
- [Ka1] W.M. Kantor: Polynomial-time algorithms for finding elements of prime order and Sylow subgroups. *J. Algorithms* **6** (1985) 478–514
- [Ka2] W.M. Kantor: Sylow's theorem in polynomial time. *J. Comp. Sys. Sci.* **30** (1985) 359–394
- [Ka3] W.M. Kantor: Finding Sylow normalizers in polynomial time. *J. Algorithms* **11** (1990) 523–563
- [Ka4] W.M. Kantor: Some topics in asymptotic group theory. *Proc. Durham Group Theory Conf.*, 1990. (To appear)
- [KL] W.M. Kantor, E. M. Luks: Computing in quotient groups. In: *Proc. 22nd ACM STOC*, Baltimore, 1990, pp. 524–534
- [Kn] D.E. Knuth: Efficient representation of perm groups. *Combinatorica* **11** (1991) 57–68 (preliminary version circulated since 1981)
- [Lub] A. Lubiw: Some *NP*-complete problems similar to graph isomorphism. *SIAM J. Comp.* **10** (1981) 11–21
- [Lu1] E.M. Luks: Isomorphism of graphs of bounded valence can be tested in polynomial time. *J. Comp. Sys. Sci.* **25** (1982) 42–65
- [Lu2] E.M. Luks: Computing the composition factors of a permutation group in polynomial time. *Combinatorica* **7** (1987) 87–99
- [LM] E.M. Luks, P. McKenzie: Fast parallel computation with permutation groups. *Proc. 27th IEEE FOCS* 1985, pp. 505–514
- [NP] P.M. Neumann, Cheryl E. Praeger: A recognition algorithm for the special linear groups. Manuscript, 1990
- [Pá] P.P. Pálfy: A polynomial bound on the orders of primitive solvable groups. *J. Algebra* **77** (1982) 127–137
- [Py] L. Pyber: The orders of doubly transitive groups: elementary estimates. (To appear)
- [Sim1] C.C. Sims: Computation with permutation groups. In: *Proc. Second Symp. Symb. Algeb. Manipulation. ACM*, New York 1971, pp. 23–28
- [Sim2] C.C. Sims: Some group theoretic algorithms. (*Lecture Notes in Mathematics* 697.) Springer, Berlin Heidelberg New York 1978, pp. 108–124
- [St] R. Steinberg: Generators for simple groups. *Can. J. Math.* **14** (1962) 277–283
- [Sto] L. Stockmeyer: The polynomial-time hierarchy. *Theoret. Comp. Sci.* **3** (1976) 1–22
- [Wi] H. Wielandt: *Finite permutation groups*. Acad. Press, New York 1964
- [Wo] T.R. Wolf: Solvable and nilpotent subgroups of  $GL(n, q^m)$ . *Can. J. Math.* **34** (1982) 1097–1111



# A Theory of Computation and Complexity over the Real Numbers

Lenore Blum<sup>1</sup>

International Computer Science Institute, 1947 Center Street, Berkeley, CA 94704;  
Mills College, Oakland, CA 94613; and University of California,  
Berkeley, CA 94720, USA

## 1. Introduction

Classically, the theories of computation and computational complexity deal with discrete problems, for example over the integers, about graphs, etc. On the other hand, most computational problems that arise in numerical analysis and scientific computation, in optimization theory and more recently in robotics and computational geometry, have as natural domains the reals  $\mathbf{R}$ , or complex numbers  $\mathbf{C}$ . A variety of ad hoc methods and models have been employed to analyze complexity issues in this realm, but unlike the classical case, a natural and invariant theory has not yet emerged. One would like to develop theoretical foundations for a theory of computational complexity for numerical analysis and scientific computation that might embody some of the naturalness and strengths of the classical theory.

Toward this goal, we have been developing a new theory of computation and complexity which attempts to integrate key ideas from the classical theory in a setting more amenable to problems defined over continuous domains. Our approach is both algebraic and concrete; the underlying space is an arbitrary commutative ring (or field) and the basic operations are polynomial (or rational) maps and tests.

The theory yields results in the continuous setting analogous to the pivotal classical results of *undecidability* and *NP-completeness* over the integers, yet reflecting the special mathematical character of the underlying space. For example, over the reals we have that (1) the Mandelbrot set as well as most Julia sets are undecidable<sup>2</sup> and (2) the problem of deciding if an algebraic variety has a real point is *NP*-complete. While there are many subtle differences between the new and classical results, the ability to employ mathematical tools of more mainstream mathematics (such as from algebra, analysis, geometry and topology) in the domain of the reals may suggest new approaches for tackling the classical, as well as new, “ $P = NP?$ ” questions.

The material covered here is based in large part on (Blum, Shub and Smale 1989) denoted in this paper by BSS, (Blum and Smale 1990) and (Blum 1990).

---

<sup>1</sup> This work was partially supported by National Science Foundation grants CCR-8712121 and CCR-8907663 and the Letts-Villard Chair at Mills College.

<sup>2</sup> Indeed, the complements of these sets provide examples of *semi-decidable* sets that are undecidable over the reals.

Discussions of related work and references are contained in those papers. See also (Shub 1990a, 1990b) and (Smale 1990). Additional relevant literature is listed in the References.

## 2. Computable Functions and Decidable Sets

The classical theory of computation had its origins in work of logicians – of Gödel, Turing, Church, Kleene, Post – in the 1930s. Of course there were no computers at the time; this work, in particular Turing’s (1937), clearly anticipated the development of the modern digital computer. But even more, a primary motivation for the logicians was to formulate and understand the concept of *decidability*, or of a *decidable* set, thus to make sense of such questions: “Is the set of theorems of arithmetic decidable?” or “Is the set of polynomials with integer coefficients and integer solutions decidable?”<sup>3</sup>

Intuitively, a set  $S \subset U$  is *decidable* if there is an “effective procedure” that given any element  $u$  of  $U$  (some natural universe) will decide in a finite number of steps whether or not  $u$  is in  $S$ , i.e. if the characteristic function of  $S$  (with respect to  $U$ ) is “effectively computable.” The models of computation designed by these logicians were intended to capture the essence of this concept of effective procedure/computation. The idea was to design formal “machines” with operations, and finitely described rules for proceeding step by step from one operation to the next, so simple and constructive that it would be self-evident that the resulting computations were effective.

In each formalism (e.g. Turing’s), a function  $f$  from the natural numbers  $\mathbb{N}$  to  $\mathbb{N}$  is defined to be *computable* if it is the *input-output* function of some such machine (e.g. a Turing machine). It is quite remarkable that even though the formalisms were often markedly different, in each case, the resulting class of computable functions (and hence decidable sets) was exactly the same. Thus, the class of computable functions appears to be a natural class, independent of any specific model of computation.<sup>4</sup> This gives one a great deal of confidence in the theoretical foundations of the theory of computation. Indeed, what is known as *Church’s thesis* is an assertion of belief that the classical formalisms completely capture our intuitive notion of computable function. Compelling motivation clearly would be required to justify yet a new paradigm.

## 3. Examples

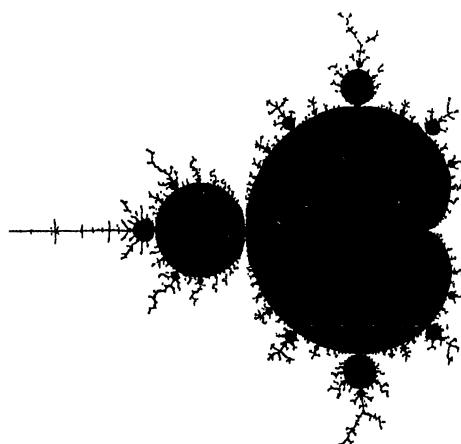
In order to motivate our theory, we briefly discuss three examples, one from complex analytic dynamics, one from numerical analysis and one from classical complexity theory.

---

<sup>3</sup> It was originally taken for granted (by mathematicians in general, and Hilbert (1901–1902) in particular) that the answers to these questions were both affirmative. The queries were actually posed as tasks: “Produce decision procedures for the given sets.” The incompleteness/undecidability results of Gödel (1931) in the first place, and of Matijasevich (1971) on the unsolvability of Hilbert’s Tenth Problem in the second, show such tasks cannot be carried out in full generality.

<sup>4</sup> These functions are often called the (*partial*) *recursive functions*.

### 3.1 Is the Mandelbrot Set Decidable?



**Fig. 1.** The Mandelbrot set<sup>5</sup>

This question was asked by Penrose (1989) in his book *The Emperor's New Mind*. Recall the Mandelbrot set  $\mathbf{M}$  can be defined:

$$\mathbf{M} = \{c \in \mathbf{C} | p_c^n(0) \nrightarrow \infty\},$$

where  $p_c(z) = z^2 + c$  and  $p_c^n$  is the  $n$ -th iterate of  $p_c$ .

It is well known (see e.g. (Branner 1989)) that the boundary of  $\mathbf{M}$  has a rich and extraordinarily complex structure. Hence, the reasonableness of Penrose's query.

However, the classical theory presupposes all underlying spaces are countable and hence *ipso facto* cannot handle such questions about arbitrary sets of real or complex numbers. One way to deal with this might be to consider the rational or algebraic skeletons of the sets in question. Problems quickly arise with this approach (e.g. consider the rational skeleton of the points on the curve  $x^3 + y^3 = 1$  in the positive orthant). Another way might be to take a recursive analysis approach. For example, we might imagine a Turing machine being input a real number bit by bit by oracle. Using its internal instructions, the machine operates on what it sees, possibly every so often outputting a bit. The resulting sequence, if any, would be considered in the limit the (binary expansion of the) real output. Problems arise here when one wants to decide if two numbers are equal.

Penrose speculates on various such approaches and concludes (p. 129) “One is left with the strong feeling that the correct viewpoint has not yet been arrived at.”

---

<sup>5</sup> This illustration is from (Penrose 1989) and is reproduced with permission by the publisher.

### 3.2 The Newton Machine

Newton's method is perhaps the "algorithm" sine qua non of numerical analysis and scientific computation. Here we briefly recall Newton's method for finding zeros of polynomials in one variable.

Given a polynomial  $f(z)$  over the complex numbers  $\mathbf{C}$ , define the *Newton map*  $N_f : S \rightarrow S$  of the Riemann sphere  $S = \mathbf{C} \cup \{\infty\}$  into itself by

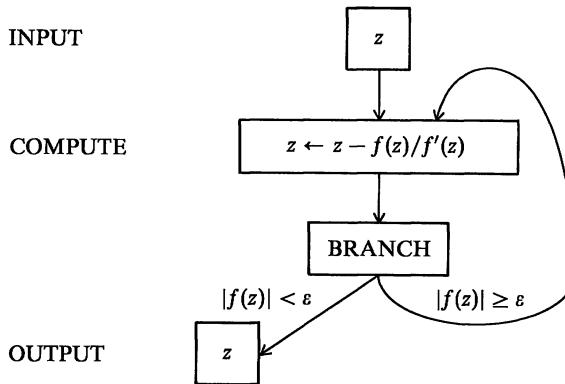
$$N_f(z) = z - (f(z)/f'(z)).$$

Now for Newton's method: Pick an initial point  $z_0 \in \mathbf{C}$  and generate the orbit

$$z_0, z_1 = N_f(z_0), z_2 = N_f(z_1), \dots, z_{k+1} = N_f(z_k) = N_f^{k+1}(z_0), \dots$$

Some stopping rule such as "stop if  $|f(z_k)| < \varepsilon$  and output  $z_k$  (else pick a new initial point if  $k$  is too large)" is implied.<sup>6</sup>

We can represent Newton's method schematically as in Fig. 2.



**Fig. 2.** The Newton machine for  $f$

A Turing machine for implementing Newton's method, by reducing all operations to bit operations, would wipe out its basic underlying structure. We would like to have a model of computation in which Newton could be represented as naturally as in the Newton machine, and in which its salient features would be as apparent.

<sup>6</sup> It is well known, however, that Newton's method is not *generally convergent*. The main obstruction to general convergence is the existence of attracting periodic points of period at least 2. (See (Smale 1985) and (Friedman 1989) for estimates on measures for the basin of attraction of the Newton map.)

### 3.3 Does $P = NP$ ?

If a problem has a solution that can be easily verified, can such a solution be found quickly? This question is formalized by means of the fundamental open problem of classical (discrete) complexity theory, namely does  $P = NP$ ? We would like to pose this question within a more general setting, thus perhaps increasing the mathematical tools and perspectives available to tackle it.

## 4. Finite Dimensional Machines over a Ring $R$

Now we describe our formal model of computation over a ring. Let  $R$  be an arbitrary ordered commutative ring (or field).

**Definition.** A finite dimensional *machine*  $M$  over  $R$  consists of three spaces: *input space*  $\bar{I}$ , *state space*  $\bar{S}$ , and *output space*  $\bar{O}$  of the form  $R^l$ ,  $R^n$ ,  $R^m$  respectively, together with a finite directed connected graph with four types of nodes: *input*, *computation*, *branch* and *output*.

The *unique* input node has no incoming edges and only one outgoing edge. All other nodes have (possibly several) incoming edges. Computation nodes have only one outgoing edge, branch nodes exactly two (left and right), and output nodes none. Each node has associated maps:

At the *input node*, there is a linear map  $I$  taking points from the input space to the state space.

Each *computation node* has an associated polynomial or rational map  $g : R^n \rightarrow R^n$  of the state space to itself.

Each *branch node* has an associated polynomial function  $h : R^n \rightarrow R$  from the state space  $R^n$  to the ring  $R$ . For a given *state*  $z$  in  $R^n$  at such a node, branching left or right will depend upon whether or not  $h(z) < 0$ .

Finally each *output node* has an associated linear map from the state space to the output space.

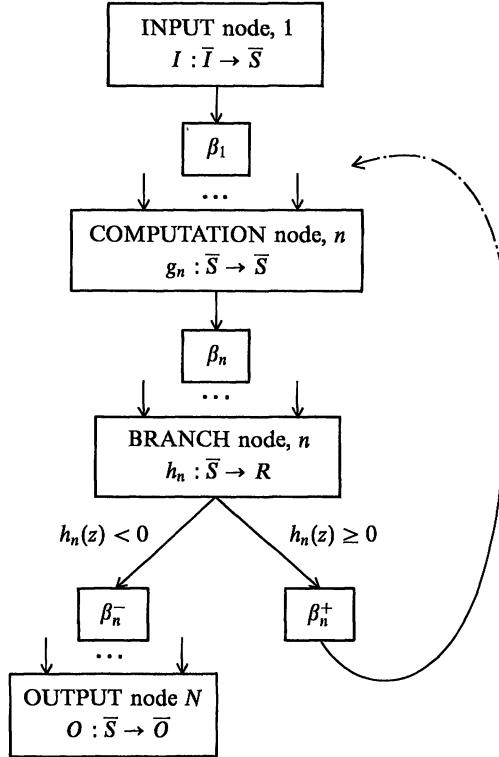
If  $R$  is a field and  $g$  is a rational map associated with some computation node, we will assume that previous nodes have tested for the vanishing of the denominators occurring in  $g$  and branched away as necessary. (Thus we are assuming that a map associated with a computation node is defined at every input to the node.)

Thus the Newton machine is an example of a machine over  $\mathbf{R}$ . (Here we are viewing  $\mathbf{C}$  as  $\mathbf{R}^2$  and the Newton map  $N_f$  as a rational map  $g = (g_1, g_2) : \mathbf{R}^2 \rightarrow \mathbf{R}^2$ . By expressing the stopping rule as  $|f(z)|^2 < \varepsilon^2$ , we get an equivalent real polynomial condition  $h(x, y) < 0$ .)

It is quite natural to view  $M$  as a discrete dynamical system. Here it is convenient to assume there is only one output node with associated map  $O$ . Thus we may let  $\bar{N} = \{1, \dots, N\}$  be the set of nodes of  $M$ , where 1 is the input node and  $N$  the output node. We call the space of node/state pairs  $\bar{N} \times \bar{S}$  the *full state space* of the machine.

Implicitly associated to  $M$  is the *computing endomorphism*

$$H : \bar{N} \times \bar{S} \rightarrow \bar{N} \times \bar{S}$$



**Fig. 3.** A finite dimensional machine  $M: I$  and  $O$  are the maps associated with the input and output nodes respectively. For  $n$  a computation node,  $g_n$  is the associated “computation map”; and for  $n$  a branch node,  $h_n$  is the associated “branching function.”

For  $n$  an input node or a computation node,  $\beta_n$  is the unique *next node* following  $n$ . For  $n$  a branch node,  $\beta_n^-$  is the next node along the left outgoing edge and  $\beta_n^+$  the next node along the right outgoing edge

of the full state space to itself. That is,  $H$  maps each node/state pair  $(n, x)$  to the unique *next node/next state* pair determined by the directed graph of the machine and its associated maps (see Fig. 3) as follows:

$$H(1, x) = (\beta_1, x); \quad H(N, x) = (N, x); \quad H(n, x) = (\beta_n, g_n(x))$$

if  $n$  is a computation node; and if  $n$  is a branch node,

$$H(n, x) = (\beta_n^-, x) \text{ if } h_n(x) < 0, \quad \text{else } (\beta_n^+, x) \text{ if } h_n(x) \geq 0.$$

The computing endomorphism is our main technical as well as conceptual tool. For example, we can use it to define the *input-output map*  $\varphi_M$  of a machine  $M$  as follows:

With *input*  $y$  in  $\bar{I}$ , let  $x = I(y)$ . Then with *initial point*  $z_0 = (1, x)$  of the full state space  $\bar{N} \times \bar{S}$  generate the *computation* (i.e. the orbit under iterates of  $H$ )

$$z_0 = (1, x), z_1 = H(z_0), z_2 = H(z_1), \dots, z_k = H(z_{k-1}) = (n_k, x_k), \dots$$

*Halt* when (if ever) the first point  $z_T$  is produced which has the form  $z_T = (N, w)$ . If this is the case, the resulting finite sequence is called a *halting computation*; we say  $M$  *halts* on input  $y$  in (*halting*) time  $T$  with *output*  $O(w)$  and define  $\varphi_M(y) = O(w)$ . If there is no such  $T$ , then  $M$  *does not halt* on input  $y$  (i.e. the halting time is infinite) and  $\varphi_M$  is not defined.

The *halting set* of  $M$ ,  $\Omega_M$ , is the set of all points in  $\bar{I}$  on which  $M$  halts. Thus,  $\varphi_M : \Omega_M \rightarrow \bar{O}$ .

The conditions describing halting computations are essentially (semi-)algebraic; they serve as the key technical tool in the proof of the *NP*-Completeness Theorem, as well as in an algebraic proof of Gödel's Theorem (see (Blum and Smale, 1990)). The basic idea is that the relevant sets can be defined in terms of these conditions. For example, the *time T halting set* of  $M$  can be defined as the set of all points  $y$  in  $\bar{I}$  for which there are solutions  $z_0, \dots, z_T$  and  $w$  to the (time  $T$ ) *register equations* (of  $M$ ):

$$z_0 = (1, I(y)), \quad z_T = (N, w) \quad \text{and} \quad z_k = H(z_{k-1}) \text{ for } k = 1, \dots, T.$$

Now having defined our formal notion of machine over  $R$ , we can easily formalize all related concepts including those in Sects. 2 and 3. For example, we define a map

$$\varphi : Y \rightarrow R^m, \quad Y \subset R^l$$

to be *computable over R* if it is the input-output map of some machine  $M$  over  $R$ , i.e. if  $\varphi = \varphi_M$  and  $Y = \Omega_M$ . We say  $M$  *computes*  $\varphi$ .<sup>7</sup> A set  $S \subset R^l$  is *decidable over R* if its characteristic function is computable over  $R$ . Otherwise it is *undecidable over R*.

In this setting, Penrose's question may thus be posed quite formally: Is the Mandelbrot set  $\mathbf{M}$  decidable over  $\mathbf{R}$ ? (Again we are viewing  $\mathbf{C}$  as  $\mathbf{R}^2$ .)

But before addressing this, it is worth noting that  $\mathbf{M}'$ , the complement of  $\mathbf{M}$ , is *semi-decidable over R*. That is, there is a machine over  $\mathbf{R}$  that on input  $x \in \mathbf{R}^2$  outputs 1 if  $x \in \mathbf{M}'$  and otherwise outputs 0 or is undefined. A semi-decidable machine for  $\mathbf{M}'$  can be constructed (see Fig. 4) using the fact that  $\mathbf{M}$  is also characterized as  $\{c \in C \mid |p_c^n(0)| \leq 2\}$ .

Now as in the classical theory, it is easy to see that a set is decidable just in case both it and its complement are semi-decidable, and that the semi-decidable sets are exactly the halting sets.<sup>8</sup>

Thus we are now ready to take a closer look at halting sets. Here it is convenient to return to the directed graph picture of a machine  $M$  over  $R$ . To each point  $y$  in the halting set  $\Omega_M$  we associate its *halting path*, i.e. the finite sequence of nodes

<sup>7</sup> We remark that the new theory reduces to the classical when  $R = \mathbf{Z}$ . That is, the computable functions over  $\mathbf{Z}$  are exactly the recursive functions (see BSS). Therefore, our model of computation is sufficiently powerful to develop the classical theory. (By Church's Thesis, we would have cause for concern had we produced more functions computable over  $\mathbf{Z}$ .)

<sup>8</sup> In the classical theory, the halting sets are exactly the *output sets* of machines. This is true also for machines over real closed fields (BSS), but not in general over arbitrary ordered rings or fields (Michaux 1990). *Open Problem.* Which results of the classical theory of computation generalize and which do not? See (Blum, Smale 1990) and (Friedman, Mansfield 1988) for additional examples.

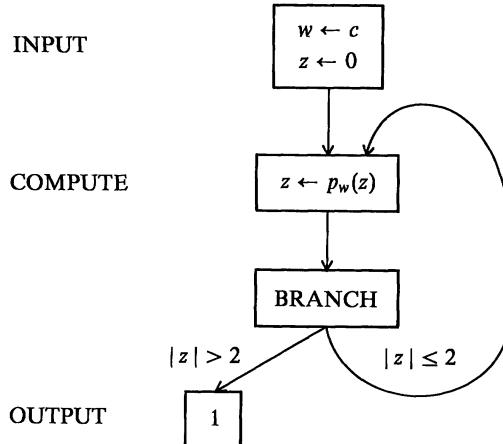


Fig. 4. A semi-decision machine for the complement of  $M$

$$n_0 = 1, n_1, \dots, n_T = N$$

traversed from input to output in the *computation* of  $\varphi_M(y)$ .

There are only a countable number of halting paths. For each halting path  $\gamma$  let  $\bar{I}\gamma$  be the set of all points in the halting set  $\Omega_M$  that have  $\gamma$  as their halting path. It is easy to see that for distinct  $\gamma$ 's the  $\bar{I}\gamma$ 's are disjoint. Also, each  $\bar{I}\gamma$  is a *semi-algebraic* set.<sup>9</sup> Note also that  $M$  acts like a “straight line program” on  $\bar{I}\gamma$ . Indeed, by concatenating the input, computation, and output maps that occur along the path  $\gamma$ , we see that  $\varphi_M$  restricted to  $\bar{I}\gamma$  is just a polynomial (or rational) map  $\varphi_\gamma$ . Thus we have the following:

**Proposition 1.** *The halting set of a machine  $M$  over  $R$  is a (disjoint) countable union of semi-algebraic sets (over  $R$ ); the input-output map  $\varphi_M$  is a piecewise polynomial (or rational) map.*

So, for example, halting sets have integral Hausdorff dimension.

**Proposition 2** (Sullivan 1990). *The Mandelbrot set is not the countable union of semi-algebraic sets over  $\mathbf{R}$ .*<sup>10</sup>

**Corollary.** *The Mandelbrot set is not decidable over  $\mathbf{R}$ .*

The same holds for most Julia sets since, from the theory of complex analytic dynamical systems, we know that most Julia sets have fractional Hausdorff dimension. Indeed, for hyperbolic rational maps of the Riemann sphere, we have the following

<sup>9</sup> A set  $S \subset R^l$  is *basic semi-algebraic* (over  $R$ ) if it is the set of elements in  $R^l$  that satisfy a (fixed) finite set of polynomial equalities and inequalities (over  $R$ ). A *semi-algebraic* set is a finite union of basic semi-algebraic sets.

<sup>10</sup> Here I would like to acknowledge helpful discussions with Michel Herman and Adrien Douady and also with John Hubbard who provided an independent proof of this result.

**Theorem (BSS).** *A Julia set is decidable if and only if it is*

1. *a round circle,*
2. *an arc of a round circle, or*
3. *the whole sphere.*

## 5. Infinite Dimensional Machines over $\mathbf{R}$

The classical construction of a universal machine assumes an effective coding of machines by (natural) numbers. In effect, the coding is a collapsing of sequences of numbers into a single number. Over the integers this can be done by a Gödel coding. However, in general over a ring  $R$  (e.g. over the reals), such an invertible collapsing cannot be done by a computable map. This suggests that a universal machine over  $R$  should have the facility to take as input finite sequences of unbounded length.

In addition, if we wish to have a natural framework for dealing with uniform procedures for solving problem instances of arbitrary dimension (say for the Travelling Salesman Problem), we are also led to consider machines that handle unbounded sequences.

With these considerations in mind we are motivated to extend our notions to *infinite dimensional machines* over  $R$ :

The underlying spaces  $\bar{I}$ ,  $\bar{S}$ , and  $\bar{O}$  for an infinite dimensional machine over  $R$  each will be  $R^\infty$ , the *infinite direct sum* space over  $R$ . A point  $y = (y_1, y_2, \dots)$  in  $R^\infty$  satisfies  $y_k = 0$  for  $k$  sufficiently large. The *length* of  $y$  is the largest  $n$  such that  $y_n \neq 0$ . Polynomial (or rational) maps in this context are still defined by a fixed finite number of polynomials (rational functions) that depend only on a fixed finite number of variables.

The machine will consist of a finite connected directed graph now containing five types of nodes, four as before, with associated maps. If the machine had only the previous type nodes it would essentially be a finite dimensional machine. The increased power comes from the addition of *fifth nodes* that allow accessing of coordinates of arbitrarily high dimension.

A *fifth node* may have several incoming edges but only one outgoing edge. The associated map transforms state  $x = (i, j, x_1, \dots, x_j, \dots, x_k, \dots, x_i, \dots)$  to state  $x' = (i, j, x_1, \dots, x_i, \dots, x_k, \dots, x_i, \dots)$ , assuming  $i$  and  $j$  are positive integers. That is, the fifth node map writes  $x_i$  in the “ $j$ -th place” of  $x$  and leaves everything else alone.

Thus, the first two coordinates of the state space play a special role which require some minor modification of the other maps. For  $y = (y_1, y_2, \dots)$  in  $\bar{I}$  we let  $I(y) = (1, 1, \text{length}(y), y_1, 0, y_2, 0, y_3, \dots)$ . This initializes the indices  $i$  and  $j$  and leaves room for workspace. Information about the length of  $y$  is often useful and so is also included. For  $x = (i, j, x_1 x_2, \dots)$  in  $\bar{S}$ , we suppose a computation node map can alter the first two coordinates only by adding 1 or by setting to 1. Finally we let  $O(x) = (x_2, x_4, \dots)$ .

Computing endomorphisms, input-output maps, halting sets and all such related notions are defined exactly as before. Note that a finite dimensional machine can be considered as a special case of infinite dimensional machine.

See BSS for an explicit construction of a universal machine over  $\mathbf{R}$ .

## 6. Complexity Theory over a Ring $R$

A goal of computational complexity theory is to quantify the intrinsic difficulty of solving problems. This theory had its origins in the 1960s.<sup>11</sup> It was developed primarily by researchers, originally trained in mathematics and logic, but who found more hospitable environments for these interests in the newly emerging computer science departments. Here, in the realm of the solvable (decidable), they discovered a rich and natural hierarchy, with the dichotomy of tractability/intractability mirroring the dichotomy of decidability/undecidability studied by the logicians. (For the seminal work in the theory, see (Rabin 1960), (M. Blum 1967), (Winograd 1970), (Cook 1971), (Karp 1972) and (Levin 1972).)

Classical complexity theory deals primarily with combinatorial (discrete, integer) problems. We extend the theory in order to consider a wider class of problems. As has been traditional however, we focus on *decision problems*. These are problems with “yes/no” answers (to questions generally of the form “Does there exist a solution to ...?”) and are classified as to their difficulty into *classes*  $P$ ,  $NP$  or as being  $NP$ -complete.

**Definition.** A *decision problem* over  $R$  is a pair  $(Y, Y_{\text{yes}})$  with

$$Y_{\text{yes}} \subset Y \subset R^\infty.$$

$Y$  is the set of *problem instances*, and  $Y_{\text{yes}}$  is the set of *yes-instances*.

For example the Travelling Salesman Problem, stated over an ordered ring  $R$ , can be put in this form by letting:

$$Y = \{(n, A, k) \mid n \text{ is a positive integer, } k > 0 \text{ and}$$

$$A = (a_{ij}) \text{ is an } n \times n \text{ matrix over } R\}$$

$$Y_{\text{yes}} = \{(n, A, k) \text{ in } Y \mid \text{there is a tour } \tau(n) \text{ with Distance } (A, \tau(n)) \leq k\}.$$

Here a tour  $\tau(n) = (\tau_1, \tau_2, \dots, \tau_n)$  is a cycle on the entire set  $\{1, 2, \dots, n\}$  and Distance  $(A, \tau(n)) = \left(\sum_{i=1}^{n-1} a_{\tau_i \tau_{i+1}}\right) + a_{\tau_n \tau_1}$ . By representing  $A$  by the sequence of its rows one after the other, we have  $Y \subset R^\infty$ .

A second example, which will be prominent in our theory, is the *4-Feasibility Problem* (*4-FEAS*) over  $R$ . Here,

$$Y = \{\text{multivariable polynomials } f \text{ over } R \mid \text{degree } f \leq 4\}$$

$$Y_{\text{yes}} = \{f \text{ in } f(\xi) = 0 \text{ for some } \xi = (\xi_1, \dots, \xi_k) \text{ in } R^k\}$$

We are supposing that polynomials are represented as elements of  $R^\infty$  via the *standard representation* (see BSS).

Thus the 4-FEAS problem is: Given a multivariable polynomial  $f$  of degree 4 with coefficients from  $R$ , does  $f(x) = 0$  have a solution over  $R$ ? While it may not at all be obvious how to decide if such a solution exists, it is a straightforward procedure to verify one that may be presented to us. Just plug the purported solution into the equation and check it out. Is this verification tractable in our model of computation? The answer will depend on the underlying mathematical

---

<sup>11</sup> However, as early as 1948, von Neumann (1963) had articulated the need for such a theory.

properties of the ring or field, as well as our measure of complexity. But first we must formalize the basic concepts of size and cost.

We first suppose we have a function *height* defined on  $R$  with values in the non-negative reals, e.g. for  $R = \mathbf{Z}$  or  $\mathbf{R}$ , and  $y \in R$  we might choose *height*( $y$ ) to be *logarithmic* height,  $\log_2(|y| + 1)$ , or *unit* height, 1. Then for  $y = (y_1, y_2, \dots, y_n, 0, 0, \dots) \in R^\infty$ , we define

$$\text{size}(y) = \text{length}(y) + \text{height}(y)$$

where  $\text{height}(y) = \max_i \text{height}(y_i)$ . Thus with unit height, *size* reflects the “dimension” of input, whereas over the integers with logarithmic height *size* reflects the traditional bit length. For the remainder of this paper, unless otherwise stated, we will suppose unit height for  $\mathbf{R}$  and logarithmic height for  $\mathbf{Z}$ .

Now suppose  $M$  is a machine over  $R$  with a height function defined. Then

$$\text{cost}_M(y) = T_M(y) \times h_{\max}(y)$$

where  $T_M(y)$  is the halting time of  $M$  on input  $y$  (which may be finite or infinite depending on whether or not  $y$  is in the halting set of  $M$ ) and  $h_{\max}(y)$  is the maximum height of any element occurring in the computation of  $M$  on input  $y$ . Over the reals, the cost reflects the number of basic algebraic operations, whereas over the integers, the cost reflects the number of bit operations.

The following definitions make sense only in case height has been defined over  $R$ .

**Definition.** A map  $\varphi$  on (*admissible*) inputs  $Y \subset R^\infty$  is *polynomial time computable* over  $R$  if there is a machine  $M$  over  $R$  that computes  $\varphi$  and

$$\text{cost}_M(y) \leq \text{poly}(\text{size}(y)), \quad \text{for all } y \text{ in } Y.$$

Here *poly* is some polynomial with nonnegative integer coefficients. Polynomial time is meant to formalize our notion of tractability.

Now we are in a position to formally define class  $P$  and class  $NP$  over  $R$ . While the first definition is straightforward, the second is considerably more subtle.

**Definition.** A decision problem  $(Y, Y_{\text{yes}})$  is in *class P (polynomial time)* over  $R$  if the characteristic function of  $Y_{\text{yes}}$  in  $Y$  is polynomial time computable over  $R$ .

**Definition.**  $(Y, Y_{\text{yes}})$  is in *class NP (non-deterministic polynomial time)* if there is a machine  $M$  which takes as input pairs  $(y, w)$  (where  $y$  in  $Y$  is a problem instance and  $w$  in  $R^\infty$  is thought of as a “guess” or “witness” for a solution to  $y$ ), outputs 1 or 0 (yes or no) and satisfies:

1. If  $y$  is a yes-instance then there exists some (guess for a solution)  $w$  such that  $\varphi_M(y, w) = 1$  and  

$$\text{cost}_M(y, w) \leq \text{poly}(\text{size}(y)).$$
2. If  $y$  is a no-instance (i.e. not a yes-instance) then there is no (guess)  $w$  such that  $\varphi_M(y, w) = 1$ .

We remark that we need only consider guesses  $w$  for which  $\text{size}(w) \leq \text{poly}(\text{size}(y))$ .

$M$  is called an *NP-decision machine* for the *NP-problem*  $(Y, Y_{\text{yes}})$ . Property 1 reflects the non-deterministic aspect of this notion, i.e. for each yes-instance, we just require that some polynomial time verifiable solution exists, not necessarily that one can be found. Property 2 requires that the verification process have some integrity, i.e. it can never output yes for a no-instance input.

In this general setting it is natural to ask (analogous to the classical question over  $\mathbf{Z}$ ): Does  $P = NP$  over  $R$ ? For  $R = \mathbf{R}$ , we have a new *open problem*.

Now let us return to 4-FEAS. The halting time for verifying a purported solution to a polynomial equation  $f(x) = 0$  using a straightforward evaluation process can easily be seen to be bounded above by a polynomial function of the length of  $f$ . (Recall we are supposing  $f$  is standardly represented as an element of  $R^\infty$ .) Thus over the reals, since size is length and cost is halting time, this verification is polynomial time.

On the other hand, over the integers we know (by the undecidability of Hilbert's Tenth Problem) that even the *smallest* size of integer solutions to polynomial equations  $f(x) = 0$  (solvable over  $\mathbf{Z}$ ) cannot be bounded above by any polynomial (in  $\text{size}(f)$ ). Thus, even if we only consider solutions of smallest size, there is no polynomial that will bound the cost of verification; in general it would just take too long to even read in purported solutions.

The above arguments show that 4-FEAS is in class *NP* over  $\mathbf{R}$  but not over  $\mathbf{Z}$ .

A key impetus for the development of classical complexity theory was the discovery (by Cook (1971) and Levin (1973)) of the existence of *NP* problems (over  $\mathbf{Z}$ ) that efficiently encode all *NP* problems.

**Definition.**  $(\widehat{Y}, \widehat{Y}_{\text{yes}})$  is *NP-complete* if it is in class *NP* and *universal* in the following sense:

For every  $(Y, Y_{\text{yes}})$  in *NP* there is a polynomial time map  $\varphi : Y \rightarrow \widehat{Y}$  such that for all  $y$  in  $Y$

$$y \text{ is in } Y_{\text{yes}} \text{ if and only if } \varphi(y) \text{ is in } \widehat{Y}_{\text{yes}}.$$

Here  $\varphi$  is the efficient (i.e. polynomial-time) *coding* function. Thus any decision procedure for  $(\widehat{Y}, \widehat{Y}_{\text{yes}})$  can be easily converted (in polynomial time) into one for  $(Y, Y_{\text{yes}})$  of not worse complexity (up to a polynomial): To decide if  $y$  is in  $Y_{\text{yes}}$ , simply encode  $y$  into  $\widehat{Y}$  using a polynomial time machine for  $\varphi$  and then decide if  $\varphi(y)$  is in  $\widehat{Y}_{\text{yes}}$ .

Thus an *NP*-complete problem is the “hardest” problem in the class *NP*; any *NP* problem can be efficiently “reduced” to it.

We have the following analogue over  $\mathbf{R}$ , to the pivotal Cook Theorem (3-SAT is *NP*-complete)<sup>12</sup> over  $\mathbf{Z}$ :

**Main Theorem (BSS).** *The 4-Feasibility Problem (4-FEAS) is *NP*-complete over the reals.*

---

<sup>12</sup>Cook's *Satisfiability Problem* is: Given a Boolean formula  $\varphi(u_1, \dots, u_k)$  is there an assignment to the variables  $u_1, \dots, u_k$  that makes the formula true? For 3-SAT the Boolean formulas considered are conjunctions of clauses of the form “ $U$  or  $V$  or  $W$ ”. Here each of  $U$ ,  $V$  or  $W$  is either a variable or the negation of a variable.

*Remarks.* This theorem has a number of immediate consequences which point to the subtle differences between the theory of  $NP$  over the integers and over the reals.

For example, over the integers it is easy to see, using a simple counting argument, that  $NP$  problems are decidable in exponential time (in the size of the instance). This is because, as noted earlier, for problem instances  $y$ , we need only consider guesses of size at most  $\text{poly}(\text{size}(y))$ . Over  $\mathbf{Z}$ , there are at most  $2^{\text{poly}(\text{size}(y))}$  such guesses, and so a perfectly good decision procedure is to check out each one in turn using an  $NP$ -decision machine for the problem.

On the other hand, over  $\mathbf{R}$  there are a continuum number of such guesses, and so it is not even clear that  $NP$  problems are decidable over  $\mathbf{R}$ , no less decidable in exponential time. However, by Tarski (1951), 4-FEAS is decidable over  $\mathbf{R}$ . So by the  $NP$ -completeness of 4-FEAS we see that all  $NP$  problems are decidable over  $\mathbf{R}$ . Moreover, (by Canny (1988) and Renegar (1988)) 4-FEAS is decidable in exponential time (over  $\mathbf{R}$ ),<sup>13</sup> and so all  $NP$  problems must be decidable in exponential time. Thus we have the same result here over the reals as over the integers but now for much deeper reasons.

The Main Theorem implies that the  $P = NP?$  problem over  $\mathbf{R}$  is equivalent to the new *open problem*: Is 4-FEAS in class  $P$  over  $\mathbf{R}$ ? (thus focusing our attention on an intrinsic algebraic-geometric problem new to complexity theory.) In contrast, recall over the integers, 4-FEAS is not *even* decidable over  $\mathbf{Z}$ .

The analogous  $NP$ -complete problem over the complex numbers  $\mathbf{C}$  is related to an effective version of Hilbert's Nullstellensatz.<sup>14</sup> Thus, as in the case of the reals, the  $NP$ -complete problem here is of a fundamental nature. However, for the moment, these are essentially the only  $NP$ -complete problems known over the reals or complex numbers.

To contrast, what makes the classical theory of  $NP$ -completeness so compelling has been the discovery (indicated first by the work of (Karp 1972)) of a large number of seemingly unrelated  $NP$ -complete problems. A polynomial time decision method for one would yield polynomial time decision methods for all.

**Open Problem.** Find other (seemingly unrelated)  $NP$ -complete problems over the reals or complex numbers.

The TSP is  $NP$ -complete over  $\mathbf{Z}$  and, as remarked earlier, in class  $NP$  over  $\mathbf{R}$ .

**Open problem:** Is TSP  $NP$ -complete over  $\mathbf{R}$ ?

**Proof of Main Theorem (Idea).** The task at hand is to show, for each  $NP$ -problem  $(Y, Y_{\text{yes}})$  over  $\mathbf{R}$ , how to encode in polynomial time any problem instance  $y$  as a degree 4 polynomial  $\hat{y}$  over  $\mathbf{R}$  such that:

$y$  is a yes-instance if and only if  $\hat{y} = 0$  has a solution over  $\mathbf{R}$ .

<sup>13</sup> For related results with respect to bit complexity see (Grigor'ev and Vorobiov 1988).

<sup>14</sup> A machine over  $\mathbf{C}$  is similar to one over  $\mathbf{R}$  except at branching nodes; over  $\mathbf{C}$  branching left or right will depend on whether or not a polynomial  $h$  evaluated at the current state  $x$  is equal to 0. The  $NP$ -complete problem over  $\mathbf{C}$  is: Given a system  $f_1, \dots, f_k$  of polynomials in  $n$  variables  $x_1, \dots, x_n$  over  $\mathbf{C}$ , decide if there is a common solution over  $\mathbf{C}$ . By Hilbert's Nullstellensatz,  $f_1, \dots, f_k$  has no common solution just in case 1 is in the ideal generated by  $f_1, \dots, f_k$ , i.e. if and only if  $1 = a_1f_1 + \dots + a_kf_k$  for some polynomials  $a_1, \dots, a_k$  over  $\mathbf{C}$  in the variables  $x_1, \dots, x_n$ .

The basic idea is to utilize the register equations for an *NP*-decision machine for  $(Y, Y_{\text{yes}})$ .

First suppose  $M$  is any machine over  $R$ . Note that the assertion “ $M$  with input  $y$  outputs  $x$  in time  $T$ ” is equivalent to asserting the system of equations

$z_0 = (1, I(y)), \quad z_T = (N, x_T), \quad O(x_T) = x \quad \text{and} \quad z_k = H(z_{k-1}) \text{ for } k = 1, \dots, T.$   
is solvable over  $R$ . We can convert (in polynomial time) this essentially semi-algebraic system over  $\mathbf{R}$  to a single polynomial equation of degree 4

$$f(y, x, u_1, \dots, u_{T'}) = 0$$

such that the original system is solvable over  $\mathbf{R}$  if and only if the single equation is. Here  $T' = p(T)$  where  $p$  is some polynomial dependent only on  $M$ . See BSS for details.

Now suppose  $M$  is an *NP*-decision machine for  $(Y, Y_{\text{yes}})$  with time bound a polynomial  $q$ . For  $y$  in  $Y$  let  $T = q(\text{size}y)$ . Suppose  $w$  is in  $\mathbf{R}^\infty$  and  $\text{size}(w) = T$ . By the above we have:

“ $M$  with input  $(y, w)$  halts with output 1 in time  $T'$ ” if and only if there is a solution to  $f((y, w), 1, u_1, \dots, u_{T'}) = 0$ . Here  $T' = p(T) = p(q(\text{size } y))$ .

Now we are ready to encode: For each  $y$  in  $Y$  let  $\hat{y}$  be the degree 4 polynomial  $f((y, w), 1, u_1, \dots, u_{T'})$  as above (having constant  $y$  and  $T + T'$  variables  $w = (w_1, \dots, w_T)$  and  $u_1, \dots, u_{T'}$ ). This is a polynomial time encoding over  $\mathbf{R}$ .

Now by the definition of *NP*-problems and *NP*-decision machines,  $y$  is in  $Y_{\text{yes}}$  if and only if: there is a  $w$  in  $\mathbf{R}^\infty$  with  $\text{size}(w) = T$  such that  $M$  with input  $(y, w)$  halts with output 1 in time  $T$ . By the above, this holds if and only if: there is a solution over  $\mathbf{R}$  to  $f((y, w), 1, u_1, \dots, u_{T'}) = 0$ , i.e. to  $\hat{y} = 0$ .

## 7. Conclusion and Directions

Since this framework is new there are a number of open problems, some of which have already been indicated and new directions to take. Many questions naturally arise concerning the relationship between the classical and new fundamental “ $P = NP?$ ” questions, and even whether the various possible extensions of the classical notion of *NP* (e.g. via guesses or via non-deterministic computation) are equivalent. (This is related to the question of whether or not the TSP is *NP*-complete over  $\mathbf{R}$ .) Here it is proving fruitful to investigate the fundamental question under varying assumptions on height, computing power and branching criteria. (See (Shub 1990b).) In the general setting, algebraic topology is providing useful tools for lower bound arguments on the topological complexity (i.e. branching complexity) of problems (Smale 1987; Vasiliev 1988; Levine 1989; Hirsch 1990). Related questions are also being pursued over other fields e.g. the  $p$ -adics (Bishop 1990).

Another direction is to study questions of parallel and distributed computation, as well as probabilistic algorithms, in this context. For the latter it would be natural to add “coin tossing” nodes to machines. Related to distributed computation, Luo and Tsitsiklis (1990) have recently given tight lower bounds for the communication complexity of several algebraic problems.

To bring the theory closer to numerical analysis and scientific computation one must extend the new model of computation to incorporate notions of round-off errors, condition numbers and approximate solutions. See (Renegar 1990)

and (Priest 1990). Here questions of the relationship between the complexity and the condition of a problem arise. In this direction it would also seem natural to adjoin nodes to compute limits of (rapidly) converging sequences, as well as other reasonable functions.

Finally there are interconnections between logic (and computation/complexity theory) and the theory of complex analytic dynamical systems to pursue. This is an intriguing direction. For example, inspired by the “degree theory” of classical recursive function theory, one is led to study the hierarchy of Julia sets imposed by various notions of relative decidability. Roughly we say a set  $A$  is *decidable relative to* a set  $B$  ( $A \leq B$ ) if a machine with an additional node for deciding  $B$  (i.e. an “oracle” for  $B$ ) can be used to decide  $A$ . The question then is: what is the resulting hierarchy? Classically, it was an open problem for a number of years (a variant of Post’s problem) to find two semi-decidable sets of integers that were incomparable with respect to relative decidability. (See (Rogers 1967).) Over  $\mathbf{R}$ , Chong (1990) has shown that the situation appears to be quite the opposite, at least for undecidable Julia sets of quadratic maps.<sup>15</sup> Thus we ask: are there two comparable undecidable Julia sets? Alternatively, is there a natural way to increase the power of machines so that the resulting hierarchy is meaningful?

In the opposite direction we have exploited the analog between computing machines and dynamical systems in our *NP*-completeness proofs over the reals (Blum, Shub and Smale 1989) and for a new proof of Gödel’s Theorem (Blum and Smale 1990). Here the computing endomorphism is our key technical tool. Can we exploit this analogy further and use techniques of dynamical systems to better understand the nature of complexity of computations and of formal computing machines? In concrete form, this approach has been successfully used by Batterson (1990) (following Shub (1983) and Smale (1985)) for the global analysis of classical algorithms of numerical linear algebra.

## References

- Batterson, S. (1990): Dynamics of eigenvalue computation. To appear in: Proceedings of the Smalefest
- Bishop, E. (1990): Computation and complexity over  $\mathbf{Q}_p$ , Work in progress
- Blum, L. (1990): Lectures on a theory of computation and complexity over the reals (or an arbitrary ring). Lectures in the Sciences of Complexity II (ed. E. Jen). Addison-Wesley, pp. 1–47
- Blum, L., Shub, M. (1986): Evaluating rational functions: Infinite precision is finite cost and tractable on average. SIAM J. Computing **15** (2) 384–398
- Blum, L., Shub, M., Smale, S. (1989): On a theory of computation and complexity over the real numbers: *NP*-completeness, recursive functions and universal machines. Bull. Amer. Math. Soc. **21**, no. 1, 1–46
- Blum, L., Smale, S. (1990): The Gödel incompleteness theorem and decidability over a ring. To appear in: Proceedings of the Smalefest
- Blum, M. (1967): A machine-independent theory of the complexity of recursive functions. J. Assoc. Comp. Mach. **14**, 322–336
- Branner, B. (1989): The Mandelbrot set. Chaos and Fractals, Proceedings of Symposia in Applied Mathematics (AMS) **39**, 75–105

<sup>15</sup> Technically, we are talking about the complements of Julia sets which are semi-decidable over  $\mathbf{R}$ .

- Canny, J. (1988): Some algebraic and geometric computations in PSPACE. Proceedings of the 20th Annual ACM Symposium on the Theory of Computing, pp. 460–467
- Chong, C.T. (1990): Reducibility of Julia sets in complex analytic dynamics. Preprint
- Cook, S.A. (1971): The complexity of theorem proving procedures. In: Proceedings 3rd ACM STOC, pp. 151–158
- Friedman, H., Ko, K. (1982): Computational complexity of real functions. *J. Theoret. Comput. Sci.* **20**, 323–352
- Friedman, H., Mansfield, R. (1988): Algorithmic procedures. Preprint, Penn State
- Friedman, J. (1989): On the convergence of Newton's method. *J. Complexity* **5**, 12–33
- Garey, M.R., Johnson, D.S. (1979): Computers and intractability. W.H. Freeman and Co., San Francisco
- Gödel, K. (1931): Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik* **38**, 173–198
- Grigor'ev, D.Yu., Vorobjov, N.N. (1988): Solving systems of polynomial inequalities in subexponential time. *J. Symb. Comput.* **5**, 37–74
- Heintz, J., Roy, M.-F., Solerno, P. (1989): Sur la complexité du principe de Tarski-Seidenberg. Preprint
- Hilbert, D. (1901–1902): Mathematical problems. *Bull. Amer. Math. Soc.* **8**, 437–479
- Hirsch, M. (1990): Applications of topology to lower bound estimates in computer science. PhD thesis, UC Berkeley
- Karp, R. (1972): Reducibility among combinatorial problems. Complexity of Computer Computations, edited by R. Miller and J. Thatcher. Plenum Press, New York pp. 85–10
- Levin, L.A. (1973): Universal sorting problems. *Problemy Peredaci Informacii* **9**, 115–116 (in Russian). [English transl.: *Problems of Information Transmission* **9**, 265–266]
- Levine, H. (1989): A lower bound for the topological complexity of poly ( $d, n$ ). *J. Complexity* **5**, 34–44
- Luo, Z.-Q., Tsitsiklis, J.N. (1990): Communication complexity of algebraic computation (extended abstract.). Proceedings of the 31st Annual Symposium on Foundations of Computer Science, IEEE Computer Society Press, pp. 758–765
- Matijasevich, Y. (1970): Enumerable sets are Diophantine. *Sov. Math. Dokl.* **11**, 354–357
- Michaux, C. (1990): Ordered rings over which output sets are recursively enumerable sets. Preprint, Université de Mons, Belgium
- Penrose, R. (1989): The Emperor's New Mind. Oxford University Press
- Pour-El, M.B., Richards, I. (1983): Computability and noncomputability in classical analysis. *Trans. Amer. Math. Soc.* **275**, 539–560
- Priest, D. (1990): Precision analysis: An approach to designing accurate computations in floating point arithmetic. Preprint
- Rabin, M.O. (1969): Degree of difficulty of computing a function and a partial ordering of the recursive sets. Technical Report No. 2. Hebrew University, Jerusalem
- Renegar, J. (1988): A faster PSPACE algorithm for deciding the existential theory of the reals. Proceedings of the 29th Annual Symposium of Computer Science. IEEE Computer Society Press, pp. 291–295
- Renegar, J. (1989): On the computational complexity and geometry of the first-order theory of the reals, Parts I, II and III. Technical Report nos. 853, 854, 856, School of Operations Research and Industrial Engineering, Cornell University, Ithaca
- Renegar, J. (1990): Some foundations for a general theory of condition numbers. To appear in: *Proceedings of the Smalefest*
- Rogers, Jr., H. (1967): Theory of Recursive Functions and Effective Computability. McGraw Hill, New York
- Shub, M. (1983): The geometry and topology of dynamical systems and algorithms for numerical problems. Lectures at Beijing University, Beijing, China

- Shub, M. (1990a): Some remarks on Bezout's theorem and complexity theory. IBM Research Report # 71048, 1990 and to appear in: Proceedings of the Smalefest
- Shub, M. (1990b): On the work of Steve Smale on the theory of computation. IBM Research Report # 71048, 1990 and to appear in: Proceedings of the Smalefest
- Smale, S. (1985): On the efficiency of algorithms of analysis. Bull. Amer. Math. Soc. **13** (2) 87–121
- Smale, S. (1987): On the topology of algorithms I. J. Complexity **3**, 81–89
- Smale, S. (1990): Some remarks on the foundations of numerical analysis. SIAM Review **32**, no. 2, 211–220
- Sullivan, D. (1990): Personal communication
- Turing, A.M. (1937): On computable numbers, with an application to the Entscheidungsproblem. Proc. Lond. Math. Soc. (Ser. 2) **42**, 230–265
- Vasiliev, V. (1988): Cohomology of the braid group and the complexity of algorithms. Preprint and to appear in: Proceedings of the Smalefest
- Von Neumann, J. (1963): The general and logical theory of automata. Collected Works, vol. V (A. Taub, ed.). MacMillan, New York, pp. 288–328
- Von zur Gathen, J. (1988): Algebraic complexity theory. Technical Report No. 207/88, Department of Computer Science, University of Toronto, 37 pp.
- Winograd, S. (1970): On the number of multiplications necessary to compute certain functions. Comm. Pure Appl. Math. **23**, 65–79



# Efficient Factoring Polynomials over Local Fields and Its Applications

Alexandre L. Chistov

Leningrad Institute for Informatics and Automation of the Academy of Sciences of the USSR (LIIAN), 14th line 39, Leningrad 199178, USSR

## Introduction

In this paper an algorithm is described for factoring multivariable polynomials over local fields. The complexity of the algorithm is polynomial in the size of input and the characteristic  $p$  of the residue field of the local field. As an application a polynomial equivalence is ascertained for the problem of constructing a basis of the ring of all integers of a given number field and the problem of finding square free part of an integer. It means that the first (respectively second) of these problems can be solved within polynomial time if there is an oracle for solving the second (respectively first) problem within polynomial time. Even the more general result is proved which is also valid in the case of non-zero characteristic, see Theorem 2.

In proofs of the last results we use on the one hand the factorization of polynomials over local fields and on the other hand an idea which is applied for obtaining efficient bounds for sizes of coefficients in the Newton-Puiseux expansion, see [8, 7] and also Lemma 1 below.

The present results solve in particular problems posed by H.W. Lenstra, Jr. in [9]. In the general case, even for one variable, earlier known algorithms required for factoring polynomials over local fields an enumeration exponential in the size of input data before applying Hensel's lemma, see [1]. Elements of local fields are represented as sums of infinite series. Here and below we regard a series as computable in time polynomial in  $A_1, \dots, A_m$  iff its  $i$ th partial sum  $S_i$  is computed in time polynomial in  $A_1, \dots, A_m$  and  $i$  for all  $i$ . Besides that, if computation of  $S_i$  involves other infinite series, then it should involve a number of initial terms polynomial in  $i$  and  $A_1, \dots, A_m$ .

Our algorithm of the factorization uses the method of Newton's polygons for constructing roots of polynomials in one variable. However, in its classical form, as in the case when the residue field is of zero characteristic, this method does not succeed because of the presence of higher ramification for extention of local fields, when one cannot choose in advance a uniformizing element in the extension. For solving the problem we use additionally expansions of a special type in the factor algebra modulo the polynomial under the factorization. Our algorithm is of the greatest interest in the case when the characteristic of the local field is zero. In the case of non-zero characteristic an analog of this algorithm is the classical algorithm

for resolution of singularities of algebraic curves over finite fields (it is also of polynomial complexity in the size of input and the characteristic of the finite field). Thus, our algorithm can be considered as a new efficient method for local resolution of singularities in rings which are finite over  $\mathbb{Z}$  of  $\mathbb{F}_p[t]$ , where  $t$  is an element algebraically independent over the field.

Now we proceed with a more detailed statement of results. Let  $k, o, \hat{k}, \delta$  and  $\pi$  denote  $\mathbb{Q}, \mathbb{Z}, \mathbb{Q}_p, \mathbb{Z}_p$  and  $p$  respectively in the case of zero characteristic, and  $\mathbb{F}_p(t)$ ,  $\mathbb{F}_p[t]$ ,  $\mathbb{F}_p((t))$ ,  $\mathbb{F}_p[[t]]$  and  $t$  in the case of non-zero characteristic (concerning standard notations see [2]). Let  $\bar{k}$  be the algebraic closure of  $\hat{k}$  and  $\text{ord} : \bar{k} \rightarrow \mathbb{Q} \cup \{\infty\}$  the order function with respect to  $\pi$ -adic metric on  $\bar{k}$ , and  $\text{ord}(\pi) = 1$ , see [1]. Let  $K'$  denote a finite separable extension of  $k$ ,  $\hat{K}$  a composite of  $K'$  and  $\hat{k}$  in  $\bar{k}$  over  $k$ ,  $\hat{O}$  the subring of  $\hat{K}$  of all the integral elements over  $\delta$ . The field  $K'$  is given over  $k$  by a primitive element  $\theta$  with its minimal polynomial  $\varphi \in k[Z]$ . Let  $\varphi_1 \in \hat{k}[Z]$  be an irreducible factor of  $\varphi$  such that  $\hat{K} \simeq \hat{k}[Z]/(\varphi_1(Z))$ , i.e.  $\varphi_1(\theta) = 0$  in  $\hat{K}$ . Without loss of generality we suppose that  $\varphi \in o[Z]$ ,  $\varphi_1 \in \delta[Z]$  and the leading coefficients  $\text{lc}_Z \varphi = \text{lc}_Z \varphi_1 = 1$ . The polynomial  $\varphi_1$  is uniquely defined by its reduction  $\bar{\varphi}_1 = \varphi_1 \bmod \pi^{r+1} \in \delta/\pi^{r+1}\delta[Z]$  where  $r = \text{ord Res}_Z(\varphi, \varphi')$  is the order of the discriminant of  $\varphi$ ; see [1], chap. 4, § 3, th. 1. The polynomial  $\varphi_1$  and consequently the field  $\hat{K}$  are given accordingly by  $\bar{\varphi}_1$ .

Furthermore (see Remark 1 below), let  $\pi_1$  be a uniformizing element of  $\hat{K}$ , and  $\eta \in \hat{K}$  such that  $\hat{k}[\eta]$  is a maximal unramified extension of  $\hat{k}$  contained in  $\hat{K}$ . Let the minimal polynomial  $h$  for  $\eta$  over  $\hat{k}$  be given, and the minimal polynomial  $g$  for  $\pi_1$  over  $\hat{k}[\eta]$ . Besides that, let  $h \in o[Z]$ ,  $\text{lc}_Z h = 1$ ,  $g \in k[\eta][Z]$ ,  $\text{lc}_Z g = 1$ . Set the field  $K = k[\eta, \pi_1]$ ,  $O = o[\eta, \pi_1]$ .

Each element  $\lambda \in \hat{K}$  can be decomposed with respect to either the  $\hat{k}$ -basis  $\{\theta^i\}_{0 \leq i < \deg \varphi_1}$ , or the  $\hat{k}$ -basis  $\{\eta^i \pi_1^j\}_{0 \leq i < \deg h, 0 \leq j < \deg g}$ , or it can be expanded into the series  $\sum_{j \geq j_{0,i}} a_{ij} \eta^i \pi_1^j$ ,  $a_{ij} \in \{0, 1, \dots, p-1\}$ . And we can go from any of these three representations to any other within polynomial time (as it was previously defined).

Let  $f \in K'[X_1, \dots, X_n]$  be an arbitrary polynomial and

$$f = \frac{1}{a} \sum_{i_1, \dots, i_n} \left( \sum_{0 \leq j < \deg \varphi} a_{i_1, \dots, i_n, j} \theta^j \right) X_1^{i_1} \dots X_n^{i_n}$$

where  $a, a_{i_1, \dots, i_n, j} \in o$ . We define the length  $l(b)$  of the element  $0 \neq b \in o$  as  $l(b) = (1 + \deg_t b)([\log_2 p] + 1)$  if  $k = \mathbb{F}_p(t)$ , and  $l(b) = \min\{s \in \mathbb{Z} : |b| < 2^{s-1}\}$  if  $k = \mathbb{Q}$ . Set  $l(0) = 1$ . If  $b \in k$ ,  $b = b_1/b_2$ ,  $\text{GCD}(b_1, b_2) = 1$  then  $l(b) = l(b_1) + l(b_2)$ . The size  $L(f)$  of  $f$  is set to be

$$l(a) + \max_{i, i_1, \dots, i_n, j} (1 + \deg_{X_i} f)^n (1 + \deg \varphi) l(a_{i_1, \dots, i_n, j})$$

Similarly we define the size  $L(\varphi)$  and the sizes of other polynomials.

Now we can formulate the main result.

**Theorem 1.** (i) In time polynomial in  $L(f)$ ,  $L(\varphi)$  and  $p$ , the decomposition  $f = \lambda \prod_{i \in I} f_i^{r_i}$  is constructed where  $f_i \in \hat{K}[X_1, \dots, X_n]$  are pairwise distinct polynomials irreducible over the field  $\hat{K}$ ,  $1 \leq r_i \in \mathbb{Z}$ ,  $I$  is a finite set, and  $0 \neq \lambda \in K$ .

(ii) When  $n = 1$  set  $X = X_1$ . Then additionally within the same time for each field  $\hat{K}_i = \hat{K}[X]/(f_i)$ ,  $i \in I$ , one can construct an element  $\eta_i \in \hat{K}_i$  with its minimal polynomial  $h_i(Z)$  over the field  $\hat{K}$  such that  $\hat{K}[\eta_i]$  is a maximal unramified extension of  $\hat{K}$  contained in  $\hat{K}_i$ , and a uniformizing element  $\pi_i$  of the field  $\hat{K}_i$  with its minimal polynomial  $g_i(Z)$  over the field  $\hat{K}[\eta_i]$ . Herewith the leading coefficients  $\text{lc}_Z h_i = \text{lc}_Z g_i = 1$  and  $h_i \in K[Z]$ ,  $g_i \in K[\eta_i][Z]$

**Remark 1.** Note that  $\varphi_1, \eta, \pi_1, h, g$  can be constructed in time polynomial in  $L(\varphi)$  and  $p$  by applying Theorem 1 to the polynomial  $\varphi \in k[Z]$  instead of  $f \in K[X]$ .

Considering derivates of polynomials Theorem 1 can be reduced to the case when  $f$  is separable, i.e. when all  $r_i = 1$ . Then for the case of one variable we obtain the following result.

**Theorem 2.** Let  $f \in \hat{O}[X]$ ,  $\varphi \in \delta[Z]$ ,  $\text{lc}_X f = \text{lc}_Z \varphi = 1$  and the order of discriminants  $\text{ord } \text{Res}_x(f, f') \leq \delta$ ,  $\text{ord } \text{Res}_Z(\varphi, \varphi') \leq \delta_1$ . Then one can construct  $f_i, \eta_i, \pi_i, g_i, h_i; i \in I$ , from Theorem 1 in time polynomial in  $\deg_X f, \deg_Z \varphi, \delta, \delta_1, p$ .

The general plan of the proof of Theorem 1 is given in Section 1.

Now we go to the applications. Let  $k, o$  be the same as above,  $f \in k[X]$  a separable polynomial,  $K' = k[X]/(f)$  the factor algebra, and  $O'$  the integral closure of  $o$  in  $K'$ . The square free part of an element  $0 \neq a \in o$  is defined as an element  $a_1 \in o$  which is equal to the product of all the irreducible divisors of  $a$  in the first power, i.e.  $a_1 = \prod_{q|a} q$  where  $q \in o$  are irreducible.

**Theorem 3.** The following problems are polynomial equivalent.

- (1) For a given element  $a \in o$  find the square free part  $a_1$  of  $a$ .
- (2) For a given separable polynomial  $f \in k[X]$  construct an  $o$ -basis of the integral closure  $O'$  of the ring  $o$  in the algebra  $K' = k[X]/(f)$ .

In the case of non-zero characteristic of  $k$  the square free part of an element  $a \in \mathbb{F}_p[t]$  can be found in polynomial time by considering the derivate  $\frac{d}{dt}$ .

Thus, the following statement is valid

**Theorem 4.** The integral closure of the ring  $\mathbb{F}_p[t]$  in a separable algebra  $\mathbb{F}_p(t)[X]/(f)$  can be constructed within the time polynomial in  $\deg_X f$  and  $\log p$ .

Note that the earlier known algorithms for constructing the integral closure from the statement of Theorem 4 required the time polynomial in  $p$  instead of  $\log p$ .

In the case when the characteristic of the field  $k$  is 0 and  $f$  is an irreducible polynomial we obtain from Theorem 4 the next

**Theorem 5.** The following problems are polynomial equivalent

- (1) For a given integer  $a$  find the square free part  $a_1$  of  $a$ .

(2) For a given irreducible polynomial  $f \in \mathbb{Q}[X]$  construct a  $\mathbb{Z}$ -basis of the ring of all integers of the field  $\mathbb{Q}[X]/(f)$ .

The polynomial reducibility of the problem (1) to the problem (2) is obtained easily by considering fields  $k(\sqrt{b})$  where  $b \in o$ . The inverse reducibility is proved much more difficult. Here we have the following

**Theorem 6.** Let  $f \in o[X]$  be a separable polynomial with the leading coefficient  $\text{lc}_X f = 1$ , and let be known the square free part  $D$  of the discriminant of the polynomial  $f$ . Then within the time polynomial in  $L(f)$ , one can construct an  $o$ -basis of the integral closure of  $o$  in the algebra  $k[X]/(f)$ .

The description of the algorithms from Theorems 3–6 is given in Section 2. The whole construction is based on the algorithm of the factorization from Theorem 1 and formulated below Lemma 1.

Let  $q \in o$  be irreducible,  $k_q$  be  $q$ -adic completion of the field  $k$ ,  $\bar{k}_q$  an algebraic closure of  $k_q$ , and  $\text{ord}_q : \bar{k}_q \rightarrow \mathbb{Q} \cup \{\infty\}$  the order function with  $\text{ord}_q(q) = 1$ . Further, let  $\text{Res}_X(f, f')$  be the discriminant of the polynomial  $f$ ,  $\text{char}(o/qo)$  the characteristic of the residue field.

**Lemma 1.** Let  $f$  be a polynomial from the statement of Theorem 6 and the degree  $\deg_X f < \text{char}(o/qo)$ . Let  $x_1, x_2 \in k_q$  be two distinct roots of the polynomial  $f$ . Then there exists an integer  $1 \leq s < \deg_X f$  for which  $\text{ord}_q \left( \frac{d^s f}{dX^s}(x_i) - b_i q^\mu \right) > \mu$ ,  $i = 1, 2$ , where  $\mu \in (1/v)\mathbb{Z}$ ,  $v \in \mathbb{Z}$ ,  $\text{GCD}(v, \text{char}(o/qo)) = 1$ ,  $1 \leq v \leq \deg_X f$ ,  $0 \leq \mu \leq (\text{ord}_q(\text{Res}_X(f, f')))/2$ ;  $b_1, b_2 \in \bar{k}_q$ ;  $\text{ord}_q b_1, \text{ord}_q b_2 \geq 0$ ;  $\text{ord}_q(b_1 - b_2) = 0$ .

An analog to this lemma (but without a restriction to  $\deg_X f$ ) is valid also in the case of zero characteristic residue fields. It was proved by the author earlier and it was central for obtaining of efficient bounds for sizes of coefficients in the Newton-Puiseux expansion, see [8]. Namely, the  $i$ th coefficient in the Newton-Puiseux expansion can be constructed within the time polynomial in the size of input and  $i^l$  where  $l$  is the transcendence degree of the field of constants over the primitive subfield. The main problem here is to estimate the first coefficients till the stabilization of the process of the expansion, i.e. the distinction of roots. Note that the direct methods give here, even in the case when there are no extensions of the constant field, an exponential or subexponential growth of sizes of coefficients. After the stabilization of the expansion one can apply Hensel's lemma in the ordinary way.

The last result has many important applications. These are algorithms with polynomial complexity for factoring multivariable polynomials over fields of power and fraction-power series and also algorithms in the theory of algebraic curves: computing the genus of a curve, indices of ramification, uniformizing elements, the smooth model of a curve, etc. within polynomial time, see [7].

## 1. Description of the Algorithm for Factoring Polynomials over Local Fields

At first we reduce the problem of factoring a polynomial  $f \in K'[X_1, \dots, X_n]$  over the field  $\hat{K}$  to the principle case when  $n = 1$ . Using the algorithm, for example from [4], one can assume  $f$  to be irreducible over  $K'$ . Applying the algorithm from [4] we can find an absolutely irreducible (i.e. irreducible over  $\hat{K}'$ ) factor  $f_1$  of  $f$ . Herewith some coefficient of  $f_1$  is equal to 1. Besides that, see [4], we construct also the field  $K'[Z]/(\psi)$  generated over  $K'$  by the coefficients of the polynomial  $f_1$ . Here  $\psi \in K'[Z]$  is an irreducible polynomial in one variable. Factor  $\psi$  over  $\hat{K}$ . Let  $\psi = \prod_{i \in I} \psi_i$  be the obtained decomposition,  $\psi_i \in \hat{K}[Z]$ , and  $\hat{K}_i = \hat{K}[X]/(\psi_i)$  the fields. Compute the norms  $f_i = N_{\hat{K}_i(X_1, \dots, X_n)/\hat{K}(X_1, \dots, X_n)}(\hat{f}_i)$  of the polynomials  $\hat{f}_i = f_1 \bmod \psi_i \in \hat{K}_i[X_1, \dots, X_n]$  for all  $i \in I$ . Then  $f = \lambda \prod_{i \in I} f_i$  is a decomposition of  $f$  into irreducible factors over the field  $\hat{K}$ , where  $0 \neq \lambda \in K'$ . Thus, everything is reduced to the factorization of  $\psi$  over  $\hat{K}$ , i.e. to the case when  $n = 1$ .

Let now  $n = 1$ , and  $f$  be a separable polynomial. Without loss of generality we can assume that  $f \in O[X]$  and  $\text{lc}_X f = 1$ . Construction of some imbedding which we are going to define is central in the proof of Theorem 1.

It is an embedding of  $\hat{K}$ -algebras

$$\hat{K}[X]/(f) \hookrightarrow \prod_{j \in J} \hat{K}_j[X]/(f_j) \quad (1)$$

where the fields  $\hat{K}_j$  are weakly ramified extensions of  $\hat{K}$ . More precisely,  $\hat{K}_j = \hat{K} \otimes_K K_j$  where  $K_j = K[\eta_j, \pi_j]$ ,  $\text{ord } \eta_j = 0$ ,  $h_j \in K[Z]$  is a minimal polynomial for the element  $\eta_j$  over  $\hat{K}$ ,  $\text{lc}_Z h_j = 1$ ,  $g_j = Z^{v(j)} - \pi_1$  is a minimal polynomial for the element  $\pi_j$  over the field  $\hat{K}[\eta_j]$ ,  $\text{GCD}(v(j), p) = 1$ . Further,  $f_j \in \hat{K}_j[X]$  is an Eisenstein's polynomial with respect to the field  $\hat{K}_j$ , and  $\deg_X f_j = p^{\varepsilon(j)}$ ,  $0 \leq \varepsilon(j) \in \mathbb{Z}$ . Therefore,  $f_j$  is irreducible over the field  $\hat{K}_j$ . Finally, for discriminants we have  $\text{ord } \text{Res}(f_j, f'_j) \leq \text{ord } \text{Res}(f, f')$  for all  $j \in J$ .  $J$  is a finite set.

To construct the embedding we need the Newton broken line of the polynomial and other notions connected with it. Now we go to their definitions. Let  $\Sigma = \{0, 1, \dots, p - 1\}$  and  $\Sigma_1 = \{\sum_{0 \leq i < \deg h} a_i \eta^i : a_i \in \Sigma\}$  be systems of representatives of the residue fields of  $\hat{k}$  and  $\hat{K}$  respectively. Let  $\psi \in O[X]$ ,  $\text{lc}_X \psi = 1$ . We write

$$\psi = \sum_{0 \leq i \in \Sigma} \sum_{0 \leq u \leq \deg \psi} a_{i,u} \pi_1^i X^u, \quad a_{i,u} \in \Sigma. \quad (2)$$

We define

$$\psi_*(N, \hat{K}) = \sum_{0 \leq i \leq N} \sum_{0 \leq u \leq \deg \psi} a_{i,u} \pi_1^i X^u \in K[X]$$

to be an approximation to the polynomial  $\psi$  in the ring  $K[X]$ ;

$$P(\psi) = \{(i, u) : a_{i,u} \neq 0\}.$$

If  $\alpha, \beta \in \mathbb{Q}$ ,  $\alpha > 0$ ,  $\beta \geq 0$  then we set  $P(\psi, \alpha, \beta) = \{(i, u) \in P(\psi) : \forall (i_1, u_1) \in P(\psi) [ai_1 + \beta u_1 \geq \alpha i + \beta u]\}$ ,

$$\psi^*(\alpha, \beta) = \sum_{(i, u) \in P(\psi, \alpha, \beta)} (a_{i,u} \bmod \pi_1) X^u \in \hat{O}/\pi_1 \hat{O}[X],$$

$$\psi(\alpha, \beta) = \sum_{(i, u) \in P(\psi, \alpha, \beta)} a_{i, u} \pi_1^i X^u \in K[X];$$

$$V(\psi) = \{(0, \deg \psi)\} \cup \{(i, u) : \exists \alpha, \beta [ \{ (i, u) \} = P(\psi, \alpha, \beta) ]\}$$

to be the set of vertices of the Newton broken line of  $\psi$ ;

$$E(\psi) = \{((i_1, u_1), (i_2, u_2)) \in V(\psi)^2 : \exists (\alpha, \beta) [ (i_1, u_1), (i_2, u_2) \in P(\psi, \alpha, \beta) \& u_1 > u_2] \}$$

to be the set of edges of the Newton broken line of  $f$ . If  $w = ((i_1, u_1), (i_2, u_2)) \in E(\psi)$ ,  $\{(i_1, u_1), (i_2, u_2)\} \subset P(\psi, \alpha, \beta)$  then we define

$$\psi_w^* = \psi^*(\alpha, \beta), \quad \lambda(w) = (u_1 - u_2)^{-1}(i_1 - i_2) = \beta/\alpha$$

to be the coefficient of the slope of the edge  $w$ , and  $0 \leq \alpha(w), \beta(w) \in \mathbb{Z}$  for which  $\beta(w)/\alpha(w) = \lambda(w)$ ,  $\text{GCD}(\beta(w), \alpha(w)) = 1$ .

Finally, if  $\psi$  is not a monomial then  $E(\psi) \neq \emptyset$ , and let  $w_0 \in E(\psi)$  be an edge for which  $\lambda(w_0) = \min\{\lambda(w) : w \in E(\psi)\}$ ,  $\psi^* = \psi_{w_0}^*$ ,  $\alpha_0 = \alpha_0(\psi, \hat{K}) = \alpha(w_0)$ ,  $\beta_0 = \beta_0(\psi, \hat{K}) = \beta(w_0)$ ,  $\alpha_1 = \alpha_1(\psi, \hat{K})$  be the greatest divisor of  $\alpha_0$  relatively prime with  $p$ .

These notions are defined (and we shall use them) for polynomials with coefficients from an arbitrary other local field  $\tilde{K}$  if  $\tilde{\Sigma}$  and  $\tilde{\pi}$  are fixed analogous to  $\Sigma_1$  and  $\pi_1$  in (2).

The embedding (1) is constructed recursively with the iteration of two cases. The first case is when  $f^*$  has two or more different roots in  $\bar{F}_p$ , the second one is when  $f^*$  is a power of a linear polynomial. In the first case we construct an intermediate embedding of  $\hat{K}$ -algebras

$$\hat{K}[X]/(f) \hookrightarrow \prod_{l \in L} \hat{K}_l[X]/(f_l) \tag{3}$$

where the fields  $\hat{K}_l$  have the same properties as  $\hat{K}_j$  above. In particular,  $\eta_l, \pi_l, h_l, g_l, \hat{O}_l = \hat{O}[\eta_l, \pi_l]$  are defined. The polynomial  $f_l^*$  is a power of a linear one for every  $l \in L$  (in the definition of  $f_l^*$  we change  $\hat{K}, \pi_1, \Sigma_1$  for  $\hat{K}_l, \pi_l, \Sigma_l = \{\sum_{0 \leq i < \deg h_l} a_i \eta_l^i : a_i \in \Sigma_1\}$  in (2)).

Now we go the construction of the embedding (3). We suppose  $f$  not to be a monomial, and besides that,  $\lambda(f, K) > 0$  changing if it is necessary  $f$  for  $\pi_1^{\deg f} f(X/\pi_1)$ . We factor  $f^*$  over  $\hat{O}/(\pi_1)$ . Namely,  $f^* = \prod_{y \in \Gamma} \bar{h}_y^{e_y}, \bar{h}_y \in \hat{O}/(\pi_1)[X]$  are irreducible,  $\text{lc}_X \bar{h}_y = 1, 1 \leq e_y \in \mathbb{Z}$ . Let  $h_y$  be a polynomial with coefficients from  $\Sigma_1$  for which  $h_y \bmod \pi_1 = \bar{h}_y, \eta_y$  a root of  $h_y$  in  $K, \pi_y^{\alpha_1(\hat{K}, f)} = \pi_1$  if  $h_y \neq X$  and  $\pi_y = \pi_1$  if  $h_y = X, K_y = K[\eta_y, \pi_y], \hat{K}_y = \hat{K} \otimes_K K_y, \pi_0^{\alpha_0(\hat{K}, f)} = \pi_1, \hat{K}'_y = \hat{K}[\eta_y, \pi_0]$ . Set  $F = f(X\pi_0^{\beta_0})/\pi_0^{\beta_0 \deg f} \in \hat{K}'_y[X]$ . Then  $F \in \hat{O}[\eta_y, \pi_0], \text{lc}_X F = 1, F^* = f^*, \lambda(F, \hat{K}'_y) = 0$ . The decomposition  $F^* = (X - \bar{\eta}_y)^{e_y} \bar{\psi}_y$  over the residue field of  $\hat{K}'_y$  is lifted by Hensel's lemma till  $F = F_y \psi_y$  over  $\hat{K}'_y$ , where  $F_y \bmod \pi_0 = (X - \bar{\eta}_y)^{e_y}$ . Set  $f_y = F_y(X/\pi_0^{\beta_0})\pi_0^{\beta_0 \deg f}$ . Then  $f_y \in \hat{K}_y[X]$  is a divisor of  $f, \text{lc}_X f_y = 1, f_y^* = (X - \bar{\eta}_y)^{e_y}$  if  $\bar{\eta}_y \neq 0$ . Compute all the polynomials  $\hat{f}_y = (f_y)_\#(\text{ord Res}(f, f')/\text{ord}(\pi_y), \hat{K}_y)$ . Then there exists an imbedding  $\hat{K}[X]/(f) \hookrightarrow \prod_{y \in \Gamma} \hat{K}_y[X]/(\hat{f}_y)$ , and  $\hat{f}_y^* = f_y^*$  if  $\bar{h}_y \neq X$ . To construct (3) it is sufficient now to do it recursively for  $\hat{f}_{y_0} \in K[X]$  with  $\bar{h}_{y_0} = X$  (if such an element  $y_0$  exists at all). Thus, we have reduced the first case to the second one.

Consider the second case which is central. Using some double recursion for each  $l \in L$  from (3) we construct a polynomial  $\tilde{f}_l \in K_l[X]$  such that  $\hat{K}_l[X]/(\tilde{f}_l) \simeq K_l[X]/(f_l)$  and either  $\tilde{f}_l^*$  is not a power of a linear polynomial, or  $\tilde{f}_l$  is an Eisenstein's polynomial with  $\deg \tilde{f}_l = p^{e(a)}$ ,  $e(l) \in \mathbb{Z}$ , and it coincides with some  $f_j$  from (2). Therefore, the second case will be reduced to the first one.

To construct  $\tilde{f}_l$  we carry out a recursive procedure within a finite number of steps beginning from step 0. At step  $a$  if it is not final we construct a polynomial  $H_{a+1} \in K_l[X]$  such that  $\text{lc}_X H_{a+1} = 1$ ,  $H_{a+1} \in \hat{O}_l[X]$ ,  $\hat{K}_l[X]/(f_l) \simeq \hat{K}_l[X]/(H_{a+1})$ ,  $\lambda(H_{a+1}, \hat{K}_l) = p^{-e(a+1)}$ ,  $0 \leq c(a+1) \in \mathbb{Z}$ ;  $c(a+1) > e(a)$  if  $a > 0$ ; and  $H_{a+1}^*$  is a power of a linear polynomial,  $\text{ord Res}(H_{a+1}, H_{a+1}') \leq \text{ord Res}(f_l, f_l')$ .

Now we describe step 0. Denote  $x = X \bmod f_l \in K_l[X]/(f_l) = \Lambda$ ,  $\lambda(f_l, K_l) = \beta_0/p^e > 0$ . We can assume without loss of generality that

$$\text{card } \Sigma_l > \deg f(\deg f + 1)/2 + 1.$$

Choose  $0 \leq v_1, v_2 \in \mathbb{Z}$ ,  $c \in \Sigma_l$ , such that  $v_2$  is minimal,  $x^{v_2}/\pi_l^{v_1} + cx$  is a primitive element of  $\Lambda$  over  $K_l$  with its minimal polynomial  $\phi$ , and  $\phi^*$  is a power of a linear polynomial,  $\lambda(\phi, \hat{K}_l) = p^{-e}$ . Set  $H_1 = \phi$ . Step 0 is not final.

We describe step  $a$  with  $a \geq 1$ . Let  $H = H_a$ ,  $\Lambda = K_l[X]/(H_a)$ ,  $x = X \bmod H_a$ ,  $H^*(\bar{\xi}) = 0$ ,  $\xi \in \Sigma_l$ ,  $\xi \bmod \pi_l = \bar{\xi}$ ,  $N_0 = p^{e(a)}$ ,  $Z_{N_0-1} = x^{N_0} \in \Lambda$ ,  $\psi_{N_0-1}$  be the minimal polynomial for  $Z_{N_0-1}$  over  $K_l$ ,  $\text{lc } \psi_{N_0-1} = 1$ . Using recursion on  $N \geq N_0 - 1$  we construct an element  $Z_N \in \Lambda$  with its minimal polynomial  $\psi_N \in K_l[X]$ ,  $\text{lc}_X \psi_N = 1$ ,  $\psi_N \in \hat{O}_l[X]$ , such that if  $N \geq N_0$  then  $Z_N = Z_{N-1} + \gamma_N x^b \pi_l^u$  where  $0 \leq b < p^{e(a)}$ ,  $b + up^{e(a)} = N$ ,  $\gamma_N \in \Sigma_l$ ,  $u, b \in \mathbb{Z}$ . Further,  $\lambda_N = \lambda(\psi_N, \hat{K}_l) > N/p^{e(a)}$ ; if  $Z_N$  is not final then  $\lambda_N = N_1/p^{e(a)}$ ,  $N_1 \in \mathbb{Z}$  and  $\psi_N^*$  is a power of a linear polynomial.

We show how to find  $\gamma_N \in \Sigma_l$ . If  $\lambda_{N-1} > N/p^{e(a)}$  then  $\gamma_N = 0$ . If  $\lambda_N = N/p^{e(a)}$  we find the root  $\bar{\xi}_{N-1}$  of  $\psi_{N-1}^*$ . Then  $\bar{\xi}_{N-1} \in \hat{O}_l/(\pi_l)$ . Set  $\bar{\gamma}_N = -\bar{\xi}_{N-1}/\bar{\xi}^b$ ,  $\gamma_N \in \Sigma_l$ ,  $\gamma_N \bmod \pi_l = \bar{\gamma}_N$ .

The process of constructing  $Z_N$  is finished at the element  $z_{N_2}$  for which one of two conditions is valid:

- (i)  $\lambda(\psi_{N_2}, \hat{K}_l) = N_3/p^{e(a)}$ , where  $e(a) > e(a)$ ,  $\text{GCD}(N_3, p) = 1$ ,  $0 \leq e(a)$ ,  $N_3 \in \mathbb{Z}$  and  $\psi_{N_2}^*$  is a power of a linear polynomial;
- (ii)  $\psi_{N_2}^*$  is not a power of a linear polynomial. We have  $N_2 < \deg f \text{ ord Res}(f, f')/(2 \text{ ord } \pi_l) = \gamma_N$ .

In both cases when (i) or (ii) is fulfilled choose  $0 \leq v_1, v_2 \in \mathbb{Z}$ ,  $c \in \Sigma_l$  such that  $v_2$  is minimal,  $z_{N_2}^{v_2}/x^{v_1} + cx$  is a primitive element of  $\Lambda$  over  $K_l$  with its minimal polynomial  $\phi \in K_l[X]$ ,  $\text{lc}_X \phi = 1$ , and  $\lambda(\phi, \hat{K}_l) = 1/\text{LCM}(p^{e(a)}, \alpha_0(\psi_{N_2}, \hat{K}_l))$ . Besides that in the case (i)  $\phi^*$  is a power of a linear polynomial and in the case (ii)  $\phi^*$  is not a power of a linear polynomial. Set  $H_{a+1} = \phi \# (\text{ord Res}(f_l, f_l')/\text{ord}(\pi_l), \tilde{K}_l)$ .

Now let condition (i) be fulfilled. Then we set  $e(a+1) = e(a)$ . Step  $a$  is not final if  $\deg H_{a+1} > p^{e(a+1)}$ . If  $\deg H_{a+1} = p^{e(a+1)}$  then step  $a$  is final, and we set  $\tilde{f}_l = H_{a+1}$ . It is Eisenstein's polynomial.

Let condition (ii) be fulfilled. Then step  $a$  is final. Set  $\tilde{f}_l = H_{a+1}$ . The polynomial  $\tilde{f}_l^*$  has at least two different roots in  $\bar{\mathbb{F}}_p$ . Thus, we have finished the description of the algorithm for constructing the embedding (1).

Having constructed (1) we can find all the roots of  $f$  contained in fields  $\hat{K}_j[X]/(f_j)$ ,  $j \in J$ , using the standard method of Newton's polygons, since we know uniformizing elements and systems of representatives of residue fields in  $\hat{K}_j[X]/(f_j)$ . After that we can find all the irreducible factors of  $f$  over  $\hat{K}$  computing, for example, normes of linear factors. Further  $\pi_i, \eta_i, h_i, g_i$  are constructed for the fields  $\hat{K}[X]/(f_i)$ ,  $i \in I$ . These fields are subfields of  $\hat{K}_j[X]/(f_j)$ ,  $j \in J$ . So here arise only some technical difficulties. We have completed the description of the algorithm from Theorem 1.

## 2. Constructing a Basis of the Ring of Integral Elements of a Global Field

For the proof of Theorems 3–6 it is sufficient to reduce problem (1) to problem (2) in Theorem 3 and to prove Theorem 6.

So let  $0 \neq b \in o$ . We find the square free part  $b_1$  of  $b$ . In the case when  $\text{char}(k) > 0$  it is possible to do this even without an algorithm for (2) considering  $\frac{d}{dt}$ . Let  $\text{char}(k) = 0$ . Set  $\theta = \sqrt{b}$ ,  $K' = k(\sqrt{b})$ ,  $\omega_1, \omega_2 \in k(\sqrt{b})$  to be  $\mathbb{Z}$ -basis of  $O'$ . We represent  $\sqrt{b} = c_1\omega_1 + c_2\omega_2$ ;  $c_1, c_2 \in \mathbb{Z}$  and compute  $b_3 = \text{GCD}(c_1, c_2)$ ,  $e_1 = c_1/b_3$ ,  $e_2 = c_2/b_3$ . Then  $\sqrt{b/b_3^2} \in O'$  and  $b_2 = b/b_3^2$  is square free. Recursively we find the square free part  $b_4$  of  $b_3$ . Then  $b_1 = \text{LCM}(b_2, b_4)$ .

Now we go to the proof of Theorem 6. We denote by  $B$  the set of all the irreducible divisors  $q \in o$  of  $D$  such that  $2 + (\deg_X f)^2 \geq \text{char}(o/qo)$  and find  $B$  by factoring  $D$  over  $\mathbb{F}_p$  or the enumeration. Set  $\delta_0 = D/\prod_{q \in B} q$ .

For each  $\delta \in o$  we consider the localizations  $O'_{(\delta)} = S^{-1}O'$  and  $o_{(\delta)} = S^{-1}o$  where  $S = o \setminus \bigcup_{q \mid \delta} qo$ . We define a  $\delta$ -integral basis of  $O'$  as a family of elements of  $O'$  which is an  $o_{(\delta)}$ -basis of  $O'_{(\delta)}$ . Consider also completions in  $\delta$ -adic topology  $o_\delta, O'_\delta, k_\delta, K'_\delta$  of  $o, O', k, K'$  respectively. Then  $\omega_1, \dots, \omega_m \in O'$  is a  $\delta$ -integral basis iff it is an  $o_\delta$ -basis of  $O'_\delta$ . Each polynomial  $0 \neq \psi \in k_\delta[X]$  can be represented in the form

$$\psi = \sum_{i_0 \leq i \in \mathbb{Z}} \sum_{0 \leq u \leq \deg \psi} a_{i,u} \delta^i X^u, \quad a_{i,u} \in \Sigma_\delta \quad (4)$$

where  $\Sigma_\delta = \{z \in \mathbb{Z} : 0 \leq z < |\delta|\}$  if  $\text{char}(k) = 0$  and  $\Sigma_\delta = \{z \in \mathbb{F}_p[t] : \deg z < \deg \delta\}$  if  $\text{char}(k) = p > 0$ ; there exists  $a_{i_0,u} \neq 0$ . Set  $\text{ord}_\delta \psi = i_0$ . Let now  $\psi \in o_\delta[X]$ ,  $\text{lc}_X \psi = 1$ . Then we define  $P(\psi), P(\psi, \alpha, \beta), \psi^*(\alpha, \beta), \psi(\alpha, \beta), V(\psi), E(\psi), \psi_w, \psi^* = \psi^*(\delta, X)$ ,  $\lambda = \lambda(\psi, \delta)$ ,  $\alpha_0 = \alpha_0(\psi, \delta)$ ,  $\beta_0 = \beta_0(\psi, \delta)$  for (4) analogous to that as it was above for (2) with changing  $\pi_1, \Sigma_1$  for  $\delta, \Sigma_\delta$ . Also we define  $\psi_\#(N, \delta)$  analogous to  $\psi_\#(N, K)$ . Further, if  $\deg \psi = 0$ ,  $\psi \in k_\delta$  we set  $\psi^* = a_{i_0,0} \bmod \delta o \in o/\delta o$

**Lemma 2.** *Let  $\psi, \psi_i \in o[X]$ ,  $\text{lc } \psi = \text{lc } \psi_i = 1$ , be separable polynomials for all  $i \in I$ . Let be given  $\omega_{i,1}, \dots, \omega_{i,m_i}$  a  $\delta$ -integral basis of the integral closure  $O_i$  of  $o$  in  $k[X]/(\psi_i)$  for every  $i \in I$ . Further, let  $\text{ord}_\delta(\psi - \prod_{i \in I} \psi_i) > (2 \deg \psi + 1) \text{ord}_\delta \text{Res}(\psi, \psi')$  where  $\text{Res}(\psi, \psi')$  is the discriminant of  $\psi$ . Then in time polynomial in  $L(\psi), L(\psi_i), L(\omega_{i,j}), 1 \leq j \leq m_i, i \in I$ , one can either construct a  $\delta$ -integral basis  $\omega_1, \dots, \omega_m$  of the integral closure  $O$  of  $o$  in  $k[X]/(\psi)$  or find a decomposition  $\delta = \delta_1 \delta_2$  where  $\delta_1, \delta_2 \in o$  are*

non-inversible in  $o$ . Besides that,  $L(\omega_i) \leq P(L(\psi), L(\delta))$ ,  $1 \leq i \leq m$ , where the polynomial  $P$  does not depend on  $\psi_i, \omega_{ij}$ .

The following lemma is an analog of Hensel's lemma.

**Lemma 3.** Let  $\psi$  be a polynomial from the statement of Lemma 2,  $\lambda = \lambda(\psi, \delta)$ ;  $g_0, h_0 \in o[X]$ ,  $\text{lc } g_0 = \text{lc } h_0 = 1$ ,  $g_0 = g_0(1, \lambda)$ ,  $h_0 = h_0(1, \lambda)$ ,  $\psi^* = g_0(1, \lambda)h_0(1, \lambda)$ ,  $\text{Res}(g_0, h_0) \bmod \delta o$  be inversible in  $o/\delta o$ . Then for every  $0 \leq i \leq m$  there exist polynomials  $g^{(i)}, h^{(i)} \in o[X]$  such that (i)  $\text{lc } g^{(i)} = \text{lc } h^{(i)} = 1$ ; (ii)  $\lambda(g^{(i)}, \delta) > \lambda, \lambda(h^{(i)}, \delta) \geq \lambda$  and the equality takes place for  $g^{(i)}$  if  $g_0$  is not a monomial and analogously for  $h^{(i)}$ ; (iii)  $g_0^*(1, \lambda) = (g^{(i)})^*(1, \lambda)$ ,  $h_0^*(1, \lambda) = (h^{(i)})^*(1, \lambda)$ ; (iv)  $\text{ord}_\delta(\psi - g^{(i)}h^{(i)}) > i/\alpha_0(\psi, \delta)$ ; (v)  $g^{(i)}, h^{(i)}$  can be found in time  $P(L(\psi), \delta, i)$ .

We need one more auxiliary algorithm which uses Euclid's algorithm for finding GCD and derivates. This algorithm has at input a polynomial  $f_0 \in o/\delta o[X]$  with  $\text{lc}_X f_0 = 1, \deg(f_0) < \text{char}(o/qo)$  for every irreducible  $q|\delta$ . It has one of three outputs:

- 1) The decomposition  $\delta = \delta_1\delta_2$  where  $\delta_1, \delta_2 \in o$  are noninversible.
- 2) The decomposition  $f_0 = f_1f_2$  where  $f_1, f_2 \in o/\delta o[X]$  are polynomials of non-zero degree  $\text{lc}_X f_1 = \text{lc}_X f_2 = 1$  and the resultant  $\text{Res}(f_1, f_2)$  is inversible in  $o/\delta o$ .
- 3) The representation  $f = f_3^\epsilon$  where  $f_3 \in o/\delta o[X], \text{lc}_X f_3 = 1, 1 \leq \epsilon \in \mathbb{Z}, f_3 \bmod q \in o/qo[X]$  are separable and  $X$  does not divide  $f_3 \bmod q$  for all  $q|\delta$ .

Now we return to the description of the algorithm from Theorem 6. We suppose without loss of generality that  $f$  is not a monomial, i.e.  $f \neq X$ . Then changing  $f$  for  $D^{\deg f}f(X/D)$  if it is necessary we assume  $\lambda(f, D) > 0$ .

For every  $q \in B$  using algorithms from section 1 we find  $f_i \in o[X]$  irreducible over  $k_q$  such that  $\text{ord}_q(f - \prod_{i \in I} f_i) > (2 \deg f + 1) \text{ord}_q \text{Res}(f, f')$ . We find also  $\eta_i^{(1)}, \pi_i^{(1)} \in k[X]/(f_i), h_i^{(1)}, g_i^{(1)}$  analogous to  $\eta_i, \pi_i, h_i, g_i$  from Section 1 with changing  $\hat{K}$  for  $k_q$ . Then  $(\eta_i^{(1)})^{m_1}(\pi_i^{(1)})^{m_2}, 0 \leq m_1 < \deg h_i^{(1)}, 0 \leq m_2 < \deg g_i^{(1)}$ , is a  $q$ -integral basis of the integral closure of  $o$  in  $k[X]/(f_i)$ ,  $i \in I$ . Applying Lemma 2 we construct a  $q$ -integral basis  $E_q$  of  $O'$  for every  $q \in B$ .

Recall that  $\delta_0 = D/\prod_{q \in B} q$ . We describe an algorithm which constructs some decomposition  $\delta_0 = \prod_{j \in J} \delta_j$  where  $\delta_j \in o$  are non-inversible and for  $j \in J$  constructs a  $\delta_j$ -integral basis  $E_j$  of  $O'$ . Then the union of families  $E_q, q \in B; E_j, j \in J; X^i \bmod f, 0 \leq i < \deg f$ , is a system of generators of the  $o$ -module  $O'$ . Using the result from [3] about the construction of a basis of a lattice over  $\mathbb{Z}$  by its system of generators in polynomial time or the similar result for the case of non-zero characteristic we find the required  $o$ -basis of  $O'$ .

Let a non-inversible divisor  $\delta$  of  $\delta_0$  be given. Now it is sufficient for us to describe an algorithm which finds either a  $\delta$ -integral basis of  $O'$  or a non-trivial decomposition  $\delta = \delta_1\delta_2$ . Indeed, if  $\delta_0$  is non-inversible then applying this algorithm to  $\delta = \delta_0$  and again to the obtained divisors  $\delta_1$  and  $\delta_2$  (if they arise) etc. we shall construct the required decomposition  $\delta = \prod_{j \in J} \delta_j$  and  $\delta_j$ -integral basises of  $O'$ .

Our algorithm is recursive on  $\deg(f)$  and under fixed degree on the maximal multiplicity of roots from  $o/qo$  of polynomials  $f^*(\delta, X) \bmod q \in o/qo[X]$  herewith the maximum is taken over all the roots and all irreducible  $q|\delta$ . In this recursive

algorithm a case may occur when  $f = X$ . Then  $f^*(\delta, X)$  is not determined but in this case  $O' = o$  and there are no difficulties.

So let  $\delta \mid \delta_0$ . At first we apply the auxiliary algorithm (see above) to the polynomial  $f^*(\delta, X)$ . If we have input 1) then the algorithm finishes its work. If we have input 2) then there exist  $g_0, h_0$  satisfying the conditions of Lemma 3 such that  $g_0^*(1, \lambda) = f_1, h_0^*(1, \lambda) = f_2, \lambda = \lambda(f, \delta)$ . We use Lemma 3 for  $\psi = f, i = (2 \deg f + 1) \operatorname{ord}_\delta \operatorname{Res}(f, f') \alpha_0(\delta, f)$  and construct  $g^{(i)}, h^{(i)}$ . We apply separately the algorithm described to  $g^{(i)}$  and  $h^{(i)}$  instead of  $f$  and further use Lemma 2.

It remains to consider input 3), i.e.  $f^*(\delta, X) = h^\varepsilon$ . At first suppose that  $\varepsilon > 1$  or  $\lambda(f, \delta)^{-1} \notin \mathbb{Z}$ . Set  $N_1 = \deg f(\deg f + 1)/2 + 1, N_2 = (\deg f)^2 + 2$  and choose subsets  $C_1, C_2 \subset \Sigma_\delta$  of  $N_1$  and  $N_2$  elements such that for them the natural through mappings  $C_i \subset \Sigma_\delta \rightarrow o/\delta o \rightarrow o/qo$  are injective for all  $q \mid \delta, i = 1, 2$ . It is possible since  $q \notin B$ . Let  $C_3 = \{y \in \mathbb{Z} : 1 \leq y \leq \deg f\}, C_4 = C_3 \times C_2 \times C_1, C'_4 = \{m\} \times C_2 \times C_1$ . We shall enumerate elements of  $C_4 \setminus C'_4$  if  $\varepsilon > 1$  and  $C'_4$  if  $\varepsilon = 1$ .

Let  $(y, c_2, c_1) \in C_4$  be an arbitrary but fixed element. We compute  $\phi_1 \in k[X]$ ,  $\operatorname{lc}_X \phi_1 = 1$  the minimal polynomial over  $k$  of the element  $f^{(y)}(x) = \frac{d^y f}{dX^y} \bmod f \in k[X]/(f) = K'$  where  $x = X \bmod f \in K'$ . We have  $\phi_1 \in o_{(\delta)}[X]$  and if  $y < m$  then  $\lambda(\phi_1, \delta) > 0$ . On the other hand  $f^{(m)}(x) = m!$  and  $\lambda(\phi_1, \delta) = 0$  for  $y = m$ .

We compute  $\mu_3 = \operatorname{LCM}(\alpha_0(f, \delta), \alpha_0(\phi_1, \delta)) > 0$  and  $v_1, v_2, v_3 \in \mathbb{Z}$  such that  $v_1 \lambda(f, \delta) + v_2 \lambda(\phi_1, \delta) + v_3 = 1/\mu_3; |v_1|, |v_2|, |v_3| \leq \mu_3 \max\{\lambda(f, \delta), \lambda(\phi_1, \delta)\}; v_3 = 0$  if  $\lambda(\phi_1, \delta) \neq 0$ . Set  $\mu_1 = \mu_3 \lambda(f, \delta) \in \mathbb{Z}, \mu_2 = \mu_3 \lambda(\phi_1, \delta)$ .

When  $y \neq m$  it is fulfilled  $\lambda(\phi_1, \delta) > 0, v_3 = 0$ . We choose  $0 < v_2 \leq \mu_1$  and define  $v_4 = v_1 - \mu_2, v_5 = v_2 + \mu_1, v_6 = 0$ .

When  $y = m$  it is fulfilled  $\lambda(\phi_1, \delta) = 0, v_3 = 0$ . We choose  $0 < v_1 \leq \mu_3, v_2 = 0$  and define  $v_4 = v_1 + \mu_3, v_5 = 0, v_6 = v_3 - \mu_1$ .

Return to the case of an arbitrary  $y$ . We compute  $\phi \in k[X], \operatorname{lc}_X \phi = 1$ , the minimal polynomial over  $k$  of the element  $z = x^{v_1}(f^{(y)}(x))^{v_2} \delta^{v_3} + c_1 x^{v_4}(f^{(y)}(x))^{v_5} \delta^{v_6} + c_2 x \in K'$ . The following lemma is based on Lemma 1 from Introduction.

**Lemma 4.** Fix an irreducible divisor  $q$  of  $\delta$ . Then there exists  $(y, c_2, c_1) \in C_4$ , among those which we enumerate, such that  $\operatorname{ord}_\delta \operatorname{Res}(\phi, \phi') \leq r = \operatorname{ord}_\delta \operatorname{Res}(f, f'), \deg \phi = m = \deg f, \lambda(\phi, q) \neq 0, \lambda(\phi, q)^{-1} \in \mathbb{Z}$  and

- (i) if  $\varepsilon > 1$  then either  $X$  divides  $\phi^*(q, X)$  or there exist more than  $\deg h = m/\varepsilon$  different roots of  $\phi^*(q, X)$  in  $\overline{o/qo}$ .
- (ii) if  $\varepsilon = 1$  then  $X$  does not divide  $\phi^*(q, X)$  and  $\phi^*(q, X)$  is separable.

Note that if  $\phi_2 = \phi^*(\delta, X) \bmod q \neq X^m$  then  $\phi^*(q, X) = \zeta^{-m} \phi_2(X\zeta), \zeta \in \overline{o/qo}$ . Lemma 4 permits either to change  $f$  for  $\tilde{f} = \phi_{\#}((2m+1)r, \delta)$  for some  $(y, c_2, c_1)$  and continue our construction recursively or obtain some decomposition  $\delta = \delta_1 \delta_2$ .

It remains to consider the case  $f = h^\varepsilon, \varepsilon = 1$  and  $\lambda(f, \delta)^{-1} \in \mathbb{Z}$ . But here a  $\delta$ -integral basis can be obtained immediately as it follows from

**Lemma 5.** Let  $\lambda(f, \delta)^{-1} = e \in \mathbb{Z}, f^*(\delta, X) \bmod q$  be separable polynomial and  $X$  does not divide  $f^*(\delta, X) \bmod q$  for all  $q \mid \delta$ . Then the family  $x^i(x^e/\delta)^j, 0 \leq i < e, 0 \leq j < m/e$ , is a  $\delta$ -integral basis of  $O'$ .

## References

1. Borevich Z.I., Shafarevich I.R.: Number theory, 2nd edn. Nauka, Moskow, 1972 (Russian) [English transl. of 1st edn.: Academic Press, 1966]
2. Bourbaki, N.: Algèbre commutative, Chap. 1–7. Actualités Sci. Indust., nos. 1290, 1293, 1308, 1314. Hermann, Paris 1961, 1964, 1965
3. Frumkin, M.A.: Application of modular arithmetics to the construction of algorithms for solving systems of linear equation. Dokl. Akad. Nauk SSSR **229** (5) (1976) 1067–1070 (Russian)
4. Chistov, A.L.: Polynomial complexity algorithm for factoring polynomials and constructing components of a variety in subexponential time. Zap. Nauchn. Semin. Leningrad. Otdel. Mat. Inst. Steklov (LOMI) **137** (1984) 124–188 (Russian) [English transl.: J. Sov. Math. **34** (4) (1986)]
5. Chistov, A.L.: Efficient factorization of polynomials over local fields. Dokl. Akad. Nauk SSSR **293** (5) (1987) 1073–1077 (Russian) [English transl.: Sov. Math. Dokl. **35** (2) (1987) 430–433]
6. Chistov, A.L.: Complexity of constructing the ring of integral elements of a global field. Dokl. Akad. Nauk SSSR **306** (5) (1989) 1063–1067 (Russian)
7. Chistov, A.L.: Polynomial complexity algorithms for computational problems in the theory of algebraic curves. Zap. Nauchn. Semin. Leningrad. Otdel. Mat. Inst. Steklov (LOMI) **176** (1989) 127–150 (Russian)
8. Chistov, A.L.: Polynomial complexity of the Newton-Puiseux algorithm. (Lecture Notes in Computer Sciences, vol. 233.) Springer, New York Berlin Heidelberg 1986, pp. 247–255
9. Open Problems from FCT'83. Preprint, Borgholm Univ., Borgholm 1983



# Interactive Proofs and Applications

Shafi Goldwasser \*

Laboratory for Computer Science, Electrical Engineering and Computer Science  
Department, Massachusetts Institute of Technology, Cambridge, MA, USA

## 1. Introduction

Proofs whose correctness can be verified efficiently play a central role in complexity theory. The famous complexity class NP consists of those sets for which “short” proofs of membership exist. For example, the set of all satisfiable Boolean formulas is in NP. A short proof that a boolean formula  $\phi$  is satisfiable would be a truth assignment to the boolean variables which makes  $\phi$  true. Formally,  $NP = \{L \subseteq \{0,1\}^* \text{ s.t. } \exists \text{ a polynomial time computable function } f_L \text{ and constant } c > 0 \text{ such that } x \in \{0,1\}^n \text{ is in } L \text{ if and only if } \exists y \in \{0,1\}^{nr} \text{ such that } f_L(x,y) = 1\}$ .

How about proving that there is *no* assignment which makes  $\phi$  true? It is generally believed that no short proof exists. Still, by some very recent work, we now know of “procedures” which can convince us quickly and beyond a shadow of a doubt that a formula  $\phi$  is not satisfiable. Such procedures, introduced by Goldwasser, Micali, and Rackoff [GMR], and in somewhat different form by Babai [Ba] are called *interactive proofs*.

Informally, an interactive proof-system is a method by which one algorithm of unlimited resources, called the *prover*, convinces another algorithm which runs in polynomial time, called the *verifier*, of the truth of a proposition. The verifier may toss coins, ask repeated questions of the prover, and run efficient tests upon the prover’s responses before deciding whether to be convinced or not. Interactive proofs do not yield proofs in the strict mathematical sense: the verifier may be incorrectly convinced with an exponentially small, though non-zero probability.

Formally, a set  $L$  is said to have an *interactive proof-system* if for all  $x$  in  $L$ , there exists a prover that can convince the verifier that  $x$  is in  $L$  with high probability, and for all  $x$  not in  $L$ , no prover can convince the verifier that  $x$  is in  $L$  with better than negligible probability. The class IP denotes the sets for which interactive proofs of membership exist.

Essentially, this procedure adds two new ingredients to the classical notion of proof (which can be written down and does not require active participation of the verifier): *randomness* and *interaction* with the prover.

---

\* Supported in part by NSF Grant 865727-CCR and ARO grant DAAL03-86-k-017. Article was partially written while on a sabbatical in Princeton University, Computer Science Department.

Lund, Fortnow, Karloff, Nisan [LFKN], and Shamir [Sh] show the exact unexpected power of adding randomness and interaction to classical proofs: a polynomial time verifier can be convinced of membership in  $L$  if and only if membership in  $L$  can be computed in polynomial space (PSPACE)<sup>2</sup>. Examples of problems that can be solved in polynomial space are showing that a boolean formula is non-satisfiable, computing the permanant of a matrix, or showing that a generalized geography game has a forced win. This particular list is of increasingly difficult “complete” problems, each as hard as a class of problems (Co-NP, FP, PSPACE).

What exactly gives interactive proofs their enhanced power? Both randomization and interaction are necessary: if the verifier did not toss coins then the prover could simply simulate all the moves of the verifier on his own and interactive proofs would be the same as the NP short proofs; in the absence of interaction, IP equals probabilistic polynomial time computation. The amount of interaction and randomness necessary is studied in [Ba, AGH, GS]. It seems that polynomial number of message exchanges are necessary to obtain the full power of interactive proofs.

Randomness is used in a different way in interactive proofs as defined in [GMR] than in the Arthur Merlin games as defined in [Ba]. In Arthur Merlin games the verifier moves are restricted to only tossing coins and sending their outcome to the prover, while in interactive proofs the verifier can toss coins in secrecy and send the prover messages computed based on the hidden coins. It is this secrecy which made possible some of the early examples of interactive proofs (see Section 4.). Still, Goldwasser and Sipser [GS] showed that the Arthur Merlin games and interactive proof systems are equal in power.

How about efficiently verifying membership in even harder sets? Ben-Or, Goldwasser, Kilian and Wigderson [BGKW] introduce procedures called *two prover interactive proofs*. Instead of one prover, it is two provers who jointly attempt to convince the verifier of the truth of a proposition. The two provers can decide on a common strategy before the interaction with the verifier starts, but once they start interacting with the verifier they can no longer communicate or see the messages exchanged between the verifier and the “other prover”. (This is reminiscent of the police practice to interrogate two suspects in a crime separately. The consistency of the alibi is what assures the police of its correctness). The class  $IP_2$  denote the sets for which membership has a two-prover interactive proof. Babai, Fortnow, and Lund [BFL] show that  $IP_2$  equals exactly nondeterministic exponential time (which contains PSPACE and is known to strictly contain NP).

The notion of interactive proof generalizes in the right way to attack a novel problem: how to convince a verifier of the truth of a proposition without giving him any extra “knowledge”. For example, convince a verifier that a Boolean formula is satisfiable without revealing a truth assignment (or *anything* he can not compute in polynomial time). This is made precise with the introduction of *zero-knowledge interactive proofs* by Goldwasser, Micali and Rackoff in [GMR]. Goldreich, Micali and Wigderson [GMW] showed that every set in NP has a zero-

---

<sup>2</sup> Note that polynomial space computation may take exponential time.

knowledge interactive proof if one way functions exist. Their result was extended [BGG, IY] to every set in  $IP$ . In contrast, [BGKW] show that zero-knowledge two-prover interactive proof exist for every set in  $IP_2$  without resorting to the unproven assumption of one-way functions.

The notion of zero-knowledge has important applications to the design of cryptographic protocols and fault tolerant computation.

## 2. Background

### 2.1 Notation

We use the symbol  $|x|$  to denote the binary length of a string  $x \in \{0, 1\}^*$ . Whenever we refer to picking an element out of a set at random we mean with uniform probability distribution. A *language*  $L$  is a subset of  $\{0, 1\}^*$ . We use the term algorithm and Turing machine interchangably throughout the paper.

### 2.2 Complexity Classes

Traditionally, efficient computation in theoretical computer science has been associated with polynomial time computation. The complexity class **P** is defined to be the set of languages for which membership can be computed by a polynomial time algorithm. (Intuitively, these languages correspond to easy to solve problems.)

In recent years probabilistic polynomial time computation has emerged as an alternative accepted formalism of efficient computation. A probabilistic algorithm is an algorithm which can toss coins as an additional primitive operation.

A language  $L$  is said to be accepted by a probabilistic polynomial time algorithm  $M$  if for all  $x \in L$ , the  $\text{prob}(M \text{ accepts } x) > \frac{2}{3}$ ; and for all  $x$  not in  $L$ , the  $\text{prob}(M \text{ rejects } x) > \frac{2}{3}$ . The class of languages accepted by probabilistic polynomial time algorithms is called **BPP**. (Intuitively, these languages correspond to problems that are easy to solve by probabilistic algorithms.)

A notable example of a problem which is in **BPP** but not known to be in **P** is integer primality testing.

The class of languages in which membership can be verified by a polynomial time algorithm is called **NP**. (Intuitively, **NP** corresponds to the set of problems for which, once solved, it is easy to verify that the solution is correct.) A language  $L \in \mathbf{NP}$  iff there exists a deterministic polynomial time algorithm  $M_L$  and a polynomial  $p_L$  such that  $x \in L$  if and only if there exists a  $y$  such that  $|y| \leq p_L(|x|)$  and  $M_L(x, y)$  accepts.

It is generally believed that  $\mathbf{P} \neq \mathbf{NP}$ . An **NP**-complete problem is a language  $L$  such that  $L \in \mathbf{NP}$  and if  $L \in \mathbf{P}$  then  $\mathbf{P} = \mathbf{NP}$ . An example of an **NP**-complete language is the set of satisfiable Boolean formulas.

The *Polynomial Time Hierarchy* (**PH**) was introduced to classify computational problems with a more complex logical structure than **NP**. It is defined inductively as follows: (in the definition below the  $|y_i| \leq \text{polynomial}(|x|)$  and  $R$  is a polynomial time computable relation.)

$$\begin{aligned}\Sigma_0^P &= \Pi_0^P = P \\ \Sigma_{k+1}^P &= \{L \mid \text{for some } R \in \Pi_k^P, L = \{x \mid \exists_y R(x, y)\}\} \\ \Pi_{k+1}^P &= \{L \mid \text{for some } R \in \sum_k^P, L = \{x \mid \forall_y R(x, y)\}\} \\ \mathbf{PH} &= \bigcup_{k>0} \sum_k^P\end{aligned}$$

A common conjecture made in complexity theory is that the hierarchy defined is strict and  $\sum_k^P \neq \sum_{k+1}^P$ .

So far, we concentrated on the time taken by an algorithm to solve a problem as the “resource” in question. Alternatively, the amount of memory used by an algorithm can be considered. We call **PSPACE** the set of languages accepted by deterministic Turing machines which are restricted to use a polynomial in the length of the input amount of space. The polynomial time hierarchy is contained in **PSPACE**, as it requires only a polynomial amount of space to go through all possible values of the universal and existential quantifiers of a logical formula to evaluate if it is satisfiable even though it requires exponential time.

### 3. Definitions

#### 3.1 Interactive Proof Systems

Before defining notion of interactive proof-systems, we define the notion of interactive Turing machine.

**Definition.** An *Interactive Turing machine (ITM)* is a probabilistic Turing machine which in addition to its input tape, random tape, work tape, and output tape, has a one read only communication tape, and one write only communication tape. The contents of the write-only communication tape can be thought of as *messages sent* by the machine; while the contents of the read-only communication tape can be thought of as *messages received* by the machine.

**Definition.** An *interactive protocol* is an ordered pair of ITMs  $(A, B)$  which share the same input tape;  $B$ 's write-only communication tape is  $A$ 's read-only communication tape and vice versa. The machines take turns in being active with  $B$  being active first. During its active stage, the machine first performs some internal computation based on the contents of its tapes, and second writes a string on its write-only communication tape. The  $i$ -th *message* of  $A(B)$  is the string  $A(B)$  writes on its write-only communication tape in  $i$ -th stage. The protocol is terminated by machine  $B$  which *accepts* (or rejects) the input by entering an accept (or reject) state. The first member of the pair,  $A$ , is a computationally unbounded Turing machine. The *computation time* of machine  $B$  is defined as the sum of  $B$ 's computation time during its active stages, and is bounded by a polynomial in the length of the input string.

A few parameters of an interactive protocol are of interest. The  $\text{value}_{A,B}(x)$  of an interactive protocol  $(A, B)$  on input  $x$  is the probability (taken over  $A$  and  $B$ 's coin tosses) that  $B$  accepts  $x$ . The number of *rounds* of an interactive protocol  $(A, B)$  on input  $x$  is defined to be the number of message exchanges between  $A$  and  $B$ . The *size* of an interactive protocol  $(A, B)$  on input  $x$  be the total number of bits exchanged between  $A$  and  $B$ .

**Definition.** We say that a language  $L \subset \{0, 1\}^*$  has an *interactive proof-system* if  $\exists \text{ITM } V \text{ s.t.}$

1.  $\exists \text{ITM } P \text{ s.t. } (P, V)$  is an interactive protocol and  $\forall x \in L$  the  $\text{value}_{(P,V)}(x) > \frac{2}{3}$ .
2.  $\forall \text{ITM } P \text{ s.t. } (P, V)$  is an interactive protocol  $\forall x \notin L$  the  $\text{value}_{(P,V)}(x) < \frac{1}{3}$ .

Note that it does not suffice to require that the verifier cannot be fooled by the predetermined prover (such a mild condition would have presupposed that the “prover” is a trusted oracle).

We say that  $(P, V)$  ( $P$  for which condition 1 holds) accepts  $L$  or is an interactive proof-system for  $L$ . Let  $\mathbf{IP}$  denote all languages accepted by some interactive proof-system. Note that  $\mathbf{NP}$  is a special case of interactive proofs, where the interaction is trivial and the verifier tosses no coins.

We say that  $(P, V)$  is a  $t(n)$ -round interactive proof-system if  $\forall x \in L$ , the number of rounds of  $(P, V)$  on input  $x$  is bounded by  $t(|x|)$ , and let  $\mathbf{IP}[t(n)]$  denote all languages accepted by some  $t(n)$ -round interactive proof-system.

We say that  $(P, V)$  has *error probability*  $\epsilon$  if  $\forall x \in L, \text{value}_{P,V}(x) > 1 - \epsilon$ , and  $\forall x \notin L, \forall P' \text{value}_{P',V}(x) < \epsilon$ .

The error probability of an interactive protocol  $(P, V)$  can be decreased to be exponentially small by a simulating<sup>3</sup> protocol  $(P', V')$  which runs the  $(P, V)$  protocol independently several times in parallel.  $V'$  accepts the input if and only if  $V$  accepts the input in a majority of the runs.

**Amplification Lemma 1.** Let  $(P, V)$  be an interactive proof-system for  $L$ . On inputs of length  $n$  to  $(P, V)$ , let  $m(n) = \text{length of messages exchanged}$ , and  $g(n) = \text{number of rounds}$ . Then  $\forall$  polynomials  $k(n)$ ,  $\exists$  interactive proof system  $(P', V')$  for  $L$  with error probability  $< \frac{1}{2^{k(n)}}$ ,  $g(n)$  rounds and length of messages exchanged  $O(m(n)k(n))$ .

### 3.2 Arthur Merlin Games (AM)

Another, seemingly more restricted, proof-system was defined by Babai [Ba1]. Babai called his proof-system an Arthur-Merlin game where Arthur corresponds to the verifier and Merlin corresponds to the prover. The difference with interactive proof-system of (GMR) is that Arthur messages consist only of the outcome of this coin tosses.

<sup>3</sup> The notion of one interactive-protocol simulating another is used throughout this paper. We say that protocol  $(P', V')$  simulates  $(P, V)$  if  $V'$  can use  $V$  as an oracle and  $P'$  can use  $P$  as an oracle and the same language is received by  $(P', V')$  as by  $(P, V)$ . The cost of the simulation is the ratio of the sizes of  $(P', V')$  and  $(P, V)$ .

An Arthur Merlin protocol (game)  $(M, A)$  on common input  $x$  is an interactive-protocol in which A's  $i$ -th message consist of the next block of  $m(|x|)$  bits on its random tape (where  $m()$  is a polynomial). The protocol is terminated by A, who then evaluates a polynomial time function defined over the contents of its tapes to decide whether to accept or reject the input  $x$ .

Thus, Arthur can be thought of as a verifier who flips public coins which Merlin, the prover, can look at.

Similarly to [GMR], Babai [Ba1] defined a language  $L$  to have an *Arthur-Merlin proof-system* if  $\exists \text{ITM } A \text{ s.t}$

1.  $\exists \text{ITM } M \text{ s.t } (M, A)$  is an Arthur Merlin Protocol and  $\forall x \in L$ , the  $\text{value}_{(M,A)}(x) > \frac{2}{3}$ .
2.  $\forall \text{ITM } M \text{ s.t } (M, A)$  is an Arthur Merlin Protocol,  $\forall x \notin L$ , the  $\text{value}_{(M,A)}(x) < \frac{1}{3}$ .

We say that an Arthur Merlin game  $(M, A)$  accepts  $L$  or is an Arthur Merlin proof-system for  $L$  if it obeys condition 1.  $\text{AM}$  will denote the set of all languages accepted by some Arthur Merlin proof-system, and  $\text{AM}[t(n)]$  denote the set of all languages accepted by a  $t(n)$ -round Arthur Merlin proof-system.

## 4. Examples of Interactive Proofs

**Notation.** Whenever an interactive protocol is demonstrated, we let  $B \rightarrow A$  : denote an active stage of machine  $B$ , in the end of which  $B$  sends  $A$  a message. Similarly,  $A \rightarrow B$  : denotes an active stage of machine  $A$ .

### Example 1: Quadratic Residuosity

Let  $Z_n^* = \{x < n, (x, n) = 1\}$

$QR = \{(x, n) \mid x < n, (x, n) \text{ and } \exists y \text{ s.t } y^2 \equiv x \pmod{n}\}$

$QNR = \{(x, n) \mid x < n, (x, n) \text{ and } \nexists y \text{ s.t } y^2 \equiv x \pmod{n}\}$

We demonstrate an interactive proof-system for  $QNR$ .

On input  $(x, n)$  to interactive protocol  $(A, B)$ :

$B \rightarrow A$  :  $B$  sends to  $A$  the list  $w_1 \cdots w_k$  where  $k = |n|$  and

$$w_i = \begin{cases} z_i^2 \pmod{n} & \text{if } b_i = 1 \\ x \cdot z_i^2 \pmod{n} & \text{if } b_i = 0 \end{cases}$$

where  $B$  selected  $z_i \in Z_n^*$ ,  $b_i \in \{0, 1\}$  at random.

$A \rightarrow B$  :  $A$  sends to  $B$  the list  $c_1 \cdots c_k$  s.t.

$$c_i = \begin{cases} 1 & \text{if } w_i \text{ is a quadratic residue mod } n \\ 0 & \text{otherwise} \end{cases}$$

$B$  accepts iff  $\forall_{1 \leq i \leq k}, c_i = b_i$

$B$  interprets  $b_i = c_i$  as evidence that  $(x, n) \in QRN$ ; while  $b_i \neq c_i$  leads him to reject.

We claim that  $(A, B)$  is an interactive proof-system for  $QNR$ . If  $(x, n) \in QNR$ , then  $w_i$  is a quadratic residue mod  $n$  iff  $b_i = 1$ . Thus, the all powerful  $A$  can easily compute whether  $w_i$  is a quadratic residue mod  $n$  or not, compute  $c_i$  correctly and make  $B$  accept with probability 1. If  $(x, n) \notin QNR$  and  $(x, n) \in QR$  then  $w_i$  is a random quadratic residue mod  $n$  regardless of whether  $b_i = 0$  or 1. Thus, the probability that  $A$  (no matter how powerful he is) can send  $c_i$  s.t  $c_i = b_i$ , is bounded by  $\frac{1}{2}$  for each  $i$  and probability that  $B$  accepts is at most  $(\frac{1}{2})^k$ .

Recognizing quadratic non-residues is in  $\textbf{NP}$ , but is not known to be solvable by probabilistic polynomial time algorithms (BPP). A short (non-interactive) proof that  $x$  is a quadratic non-residue modulo  $n$  is a prime factor  $p$  of  $n$  such that the Legendre symbol of  $x \bmod p$  is  $-1$ .

The interest of the interactive proof for  $QNR$  described above is that none of the prime factors of  $n$  are disclosed to the verifier during the interactive proof. In fact, this example was the first zero-knowledge interactive proof (see [GMR] for definition) known for a language not known to be solvable in polynomial time. This was shown by Goldwasser Micali and Rackoff in [GMR].

### Example 2: Graph Non-Isomorphism

The most famous interactive proof for a problem not known to be in  $\textbf{NP}$  is for the *graph non-isomorphism* problem. This was shown by Goldreich, Micali, and Wigderson [GMW]. The input is a pair of graphs  $G_1$  and  $G_2$ , and one is required to prove that there exists no 1-1 edge-invariant mapping of the vertices of the first graph to the vertices of the second graph. (A mapping  $\pi$  from the vertices of  $G_1$  to the vertices  $G_2$  is *edge-invariant* if the nodes  $v$  and  $u$  are adjacent in  $G_1$  iff the nodes  $\pi(v)$  and  $\pi(u)$  are adjacent in  $G_2$ .)

The length of the shortest (non-interactive) proof of non-isomorphism is no better than the best deterministic isomorphism test, i.e  $e^{O(\sqrt{n \log n})}$ .

The interactive proof  $(A, B)$  on input  $(G_1, G_2)$  proceeds as follows:

$B \longrightarrow A : B$  chooses at random one of the two input graphs,  $G_{\alpha_i}$

where  $\alpha_i \in \{1, 2\}$ .  $B$  creates a random isomorphic copy of  $G_{\alpha_i}$  and sends it to  $A$ . (This is repeated for  $1 \leq i \leq k$ , with independent random choices, where  $k =$  number of vertices in  $G_{\alpha_i}$ )

$A \longrightarrow B : A$  sends  $B \beta_i \in \{1, 2\}$  for all  $1 \leq i \leq k$ .

$B$  accepts iff  $\beta_i = \alpha_i$  for all  $1 \leq i \leq k$ .

$B$  interprets  $\beta_i = \alpha_i$  as evidence that the graphs are not isomorphic; while  $\beta_i \neq \alpha_i$  leads him to reject.

If the two graphs are not isomorphic, the prover has no difficulty to always answer correctly (i.e., a  $\beta$  equal to  $\alpha$ ), and the verifier will accept. If the two graphs are isomorphic, it is impossible to distinguish a random isomorphic copy of the first from a random isomorphic copy of the second, and the probability that the prover answers correctly to one “query” is at most  $\frac{1}{2}$ . The probability that the prover answers correctly all  $k$  queries is  $\leq (\frac{1}{2})^k$ .

## 5. Equivalence of Interactive Proof Systems and Arthur Merlin Proof Systems

If we review the examples of interactive-proof systems for quadratic non-residuosity and graph non-isomorphism, we see that the ability of the verifier to keep his coin flips secret, and thus unable the prover to predict the correct answer to his questions unless the statement in question was true, was the *crucial* ingredient which made these proofs go through. Thus, the interactive proofs described in the example are not Arthur Merlin games. This leads to our next question: as Arthur is a special case of a general verifier who can not hide the results of his coin flips, **AM** is contained in **IP**, but is it a strict containment?

Surprisingly, it turns out that **AM** = **IP**. Goldwasser and Sipser [GS] show a transformation from a general interactive proof-system to an Arthur Merlin protocol accepting the same language. Moreover, the transformation preserves the number of rounds.

**Theorem 2** [GS].  $\forall$  polynomially bounded  $t(n) \geq 2$   $AM[t(n)] = IP[t(n)]$ .

In particular, the theorem implies that the graph non-isomorphism language is accepted by a constant round Arthur Merlin protocol. This has certain implication as to the complexity of graph non-isomorphism (see Section 7.2 ).

The equivalence between **AM** and **IP** is highly convenient. It is usually easier to prove membership of languages using the interactive proof formulation. The elegant simplicity of the Arthur-Merlin games definition facilitates proving complexity results.

For example, we have defined **IP** to allow two sided error: the verifier can err both when  $x$  is in the language and when  $x$  is not in the language. Goldreich, Sipser and Mansour [GSM] show that every  $L \in AM[t(n)]$  has a one sided error Arthur Merlin protocol  $(M, A)$  of  $t(n)$ -rounds. Namely, when  $x \in L$ , the  $\Pr(A \text{ accepts } x) = 1$ .

## 6. Interaction: Essential Ingredient?

When the verifier is deterministic it is clear that interaction does not add power with respect to language recognition, since the all powerful prover can anticipate in advance all messages of the verifier and answer them with no need for interaction. However, when the verifier is a probabilistic algorithm, it is an interesting and not fully resolved question whether more rounds of interaction enable a proof-system to recognize more languages. It has been shown by Babai [Ba1], that any language which has a constant round interactive proof-system, has an interactive proof-system that requires only two rounds.

**Theorem 3** [Ba]. For constant  $c \geq 2$   $AM[c] = AM[2]$ .

Babai and Moran [BM] obtained a stronger result showing how the number of rounds can be halved .

**Theorem 4** [BM]. *For polynomially bounded  $t(n) \geq 2$   $\text{AM}[2t(n)] = \text{AM}[t(n) + 1]$ .*

The simulation technique used to prove the collapse and speed-up theorem result in a  $O(t(n))$  cost. Thus, if the number of rounds is halved a logarithmic number of times, the size of game grows faster than a polynomial in  $n$ . Thus, the speed-up theorem can not resolve the question whether for unbounded  $t(n)$ ,  $\text{AM}[t(n)] = \text{AM}[2]$ .

One piece of evidence known in regard to this question was obtained by Aiello, Goldwasser, and Hastad [AGH] as follows.

**Theorem 5** [AGH].

$$\forall f, g \text{ s.t } g(n) = o(f(n)), \exists \text{ oracle set } B \text{ s.t } \text{AM}^B[f(n)] \neq \text{AM}^B[g(n)]$$

This result does not imply any unrelativized conclusion, but implies that resolving the above question will involve proof techniques which do not relativize (e.g., diagonalization type arguments will not suffice). As a word of warning we note that quite recently a few relativized results in this field have shown not to hold in a non-relativized setting. See Section 8..

## 7. The Complexity of Interactive Proofs

An interesting question is what is the complexity of language in **IP** in terms of more traditional complexity classes. For example, can a verifier be convinced that a Boolean formula  $\phi$  has *no* satisfying assignments, or that the *shortest* traveling salesman tour in a graph is bounded by integer  $k$ ?

The answer to this question is different for interactive proofs with bounded (constant independent of the input) number of rounds and interactive proofs with polynomial number of rounds.

We start with the case of bounded rounds.

### 7.1 Bounded Rounds Interactive Proofs

The *Polynomial Time Hierarchy* (PH) (see Section 2.2) was introduced to classify computational problems with a more complex logical structure than **NP**. We can extend the definition of  $\sum_k^P$  and  $\prod_k^P$  from a finite number of alternation of quantifiers, to a number of alternations which is a function of the input length. Namely,  $\sum_{t(n)}^P = \{L | L = \{x | \exists y_{t(n)} \forall y_{t(n)-1} \dots R(x, y_1 \dots, y_{t(n)})\}, R \in P\}$  and  $\prod_{t(n)}^P = \{L | L = \{x | \forall_{y_{t(n)}} \exists_{y_{t(n)-1}} \dots R(x, y_1 \dots, y_{t(n)})\}, R \in P\}$ .

It can be shown (using ideas from a proof of Lautman [L]) that membership in any language accepted by an interactive proof of  $t(n)$  rounds can be expressed as the question of satisfying a logical formula with  $t(n)$  alternation of universal and existential quantifiers starting with a universal quantifiers.

**Theorem 6** [Ba, AGH].  $\text{AM}[t(n)] \subseteq \prod_{t(n)}^P$  for all  $t(n) \geq 2$ .

This was shown by Babai [Ba1] for  $t(n) = 2$ , and Aielo-Goldwasser-Hastad for unbounded  $t(n)$ . The converse is not known to be true.

An especially interesting case in question is whether languages in  $\text{Co-NP} = \prod_1^P$  are in  $\text{IP}[2]$ . A result which may be taken as supporting a separation between  $\text{Co-NP}$  and  $\text{IP}[2]$  is by Boppana-Hastad-Zachos [BHZ].

**Theorem 7** [BHZ]. If  $\text{Co-NP} \subseteq \text{AM}[2] \implies \text{PH} = \prod_2^P$ .

We conclude that showing that  $\text{Co-NP}$  has bounded round interactive proofs would be quite hard if true. The proof of Theorem 7, however, does not carry over to  $\text{AM}[t(n)]$  for the case of unbounded number of rounds  $t(n)$ .

A language is said to have non-uniform polynomial time deterministic (non-deterministic) algorithms if for every size of input  $n$ , there exists a polynomial in  $n$  time deterministic (non-deterministic) algorithm which decides membership of strings of length  $n$  in the language. Non-uniform complexity has generated much interest in complexity theory in the last years.

**Theorem 8.** *Languages in  $\text{IP}[2]$  have non-uniform polynomial-time non-deterministic algorithms.*

## 7.2 The Complexity of Graph Non-Isomorphism

The graph isomorphism problem is one of a select group of problems in  $\text{NP}$  which have not been classified as solvable in polynomial time, nor have been proven  $\text{NP}$ -complete.

The developments regarding interactive proof-system have shed light on the complexity of the graph non-isomorphism problem.

Let  $\text{ISO} = \{\text{graphs } (G_1, G_2) \text{ s.t } G_1 \text{ is isomorphic to } G_2\}$  and  $\text{Non-ISO}$  as its complement.

In Section 4. we showed that  $\text{Non-ISO}$  has a bounded round interactive proof and thus is in  $\text{IP}[2]$ . Combining the results above we can show the following conditional implications.

**Corollary 9.** *If  $\text{ISO}$  is  $\text{NP}$ -complete then*

1. *Polynomial Time Hierarchy ( $\text{PH}$ ) collapses to  $\prod_2^P$ .*
2.  *$\text{Co-NP}$  has non-uniform polynomial time non-deterministic algorithms.*

This is usually taken as evidence that  $\text{ISO}$  is *not*  $\text{NP}$ -complete. In general, showing that both a language  $L$  and its complement  $\bar{L}$  are in  $\text{IP}[2]$  imply that either  $L$  is not  $\text{NP}$ -complete or the polynomial time hierarchy collapses and  $\text{Co-NP}$  has non-uniform polynomial size non-deterministic algorithms. This can help when trying to classify the complexity of an  $\text{NP}$  problem. An interesting  $\text{NP}$  problem which is not known to be  $\text{NP}$ -complete is the problem of the existence of a short vector in an integer lattice. It is open question whether the complement problem (all vectors in a lattice are bigger than a given value) is in  $\text{IP}[2]$ .

## 8. IP = PSPACE

We saw in the previous section that bounded round interactive proofs are quite low in the polynomial time hierarchy. In this section, we investigate how powerful are general interactive proofs.

We first bound the complexity of languages accepted by interactive proofs. Membership in every  $L \in \text{IP}$  can be decided by some Turing machine which uses a polynomial amount of space and  $\text{IP} \subseteq \text{PSPACE}$  as follows.

Let  $L$  be in  $\text{IP}$  accepted by an interactive proof  $(P, V)$  with error probability  $\epsilon$ . On input  $x$  a Turing machine can in polynomial space go through all possible histories of communication between  $P$  and  $V$  and compute the probability that the verifier accepts. If this probability is greater than  $\epsilon$ , then  $x$  is in  $L$ .

Until recently, it was not known whether the containment of  $\text{IP} \subseteq \text{PSPACE}$  is strict. In fact, only a handful of languages not known to be in  $\text{NP}$  were shown to be in  $\text{IP}$ .

In December 1989, Lund Fortnow Karloff and Nisan showed that every language in the polynomial time hierarchy has an interactive proof system. Their proof used novel techniques showing how to reduce proving to the verifier the number of accepting computations of a non-deterministic Turing machine, to proving the value of polynomials of low degree over a finite field. These algebraic techniques are of general interest to complexity theory at large as they do not relativize.

**Theorem 10** [LFKN].  $PH \subseteq IP$ .

The result was announced to a few colleagues using the international electronic mail system. A few days later A. Shamir [Sh] closed the gap and showed that in fact languages accepted by interactive proofs with polynomial number of rounds of interaction are exactly those languages accepted in polynomial space.

**Theorem 11** [Sh].  $IP = PSPACE$ .

Shamir's proof shows an interactive proof for the *quantified Boolean formula* (QBF) language defined as follows:

$$QBF = \{F = (Q_1 y_1, Q_2 y_2, Q_n y_n, R(y_1, \dots, y_n))\},$$

where  $R$  is an unquantified Boolean formula and  $Q_i \in \{\exists, \forall\}$ . Membership in any language in  $\text{PSPACE}$  can be reduced to membership in QBF.

We briefly sketch the ideas in the proof. The idea of [LFKN] can be described as follows. Let  $L \in \Pi_k^P$ . Then,  $x \in L$  if and only if  $F(x) = \forall y_1 \exists y_2 \dots \exists y_k R(x, y_1, \dots, y_k)$  is true. (Here  $k$  is *finite*, and  $R$  an unquantified Boolean formula). We can write an arithmetic expression  $\tilde{F} = \prod_{y_1=0}^1 \sum_{y_2=0}^1 \dots \sum_{y_k=0}^1 \tilde{R}(x, \tilde{y}_1, \dots, \tilde{y}_k)$ , replacing every Boolean variable  $y_i$  by integer variable  $\tilde{y}_i$ , every logical  $\wedge$  to arithmetic  $\times$ , every logical  $\neg E$  to  $1 - E$ , every  $\forall y$  to  $\prod_{y \in \{0,1\}}$ , and every  $\exists y$  to  $\sum_{y \in \{0,1\}}$ .  $\tilde{F}$  evaluates to 0 if and only if the input  $x$  does not satisfy  $F$ . Thus in order to prove that  $x$  satisfies  $F$  the prover has to prove that  $\tilde{F} \neq 0$ . Define  $g_{y_1}(z) = \tilde{F}_{y_1=z}$

(evaluated  $\tilde{F}$  with  $y_1$  replaced everywhere by variable  $z$ ). Note that since  $k$  is finite,  $g_{y_1}(z)$  is a polynomial of polynomial degree in variable  $z$ . The prover sends the coefficients of  $g_{y_1}(z)$  to the verifier. The verifier checks that the value of  $g$  received on  $z = 0, 1$  matches the assertion that  $\tilde{F} \neq 0$ . If this assertion was false, then it must be that the  $g$  received is not the correct polynomial, and since both  $g$  received and the correct polynomial are of low degree they can agree on only few points. Now, the verifier chooses a random value  $z'$  from some sufficiently large range of integers, and asks the prover to prove that the value of  $g_{y_1}(z')$  equals  $\tilde{F}_{y_1/z'} (\tilde{F} \text{ when } y_1 \text{ is assigned value } z')$ . Since  $z'$  was chosen randomly among integers from a sufficiently large range, with high probability this new assertion is false too. The verifier and prover iterate on that until they exhaust all the variables  $y_i$ , and are left with a fully instantiated expression, at which point the verifier can verify whether the value the prover claims for it is true.

When extending this proof to **PSPACE**, the number of quantifiers  $k$  is not finite, and the degree of the polynomial  $\tilde{F}$  and the intermediate  $g$ 's in the variables  $y_i$  can be exponential! Shamir's [Sh] solution is to replace expressions of the form  $(\forall y_i R)$  by  $(\forall y_i \exists y'_1, \dots, y'_{i-1}, \text{ such that } (y_1 = y'_1, y_2 = y'_2, \dots, y_{i-1} = y'_{i-1}, y_i = y'_i) \wedge R_{y=y})$  (i.e replace each  $y_l$  in  $R$  by  $y'_l$ ,  $l \leq i$  in  $R$ ). This way, no variable in  $\tilde{F}$  will have higher than constant degree and the number of variables is squared. When  $k$  is not finite the value of  $\tilde{F}$  may be doubly exponential in size, thus the prover will prove to the verifier that  $\tilde{F} \neq 0$  modulo a prime of smaller size. The above procedure is modified accordingly.

A few remarks are in order. First, the equivalence with **PSPACE** implies that **IP** is closed under complementation. Second, a polynomial number of rounds of interaction seem to be crucial for the argument used in the proofs of Theorems 10 and 11. Showing that **IP** = **IP**[2] would imply that **PSPACE** =  $\Pi_2^P$ . We conclude that resolving the question of whether bounded round interactive proofs accept all of **IP** would settle questions complexity theorists have been struggling with for years.

The current proof technique of Theorem 11 requires the prover to be **PSPACE** powerful. Possible conjectures are (1) that the ability to decide whether or not a Boolean formula is satisfiable does not enable a prover to convince a polynomial time verifier that a formula is not satisfiable (i.e membership in **Co-NP** complete languages), and (2) that the ability to *count* the number of satisfying assignments to a Boolean formula is sufficient to enable a prover to convince a verifier of membership in **PSPACE** complete languages. It seems likely progress can be made on these questions in near future.

## 9. Multi-Prover Interactive Proofs

Is membership in a **PSPACE** language the limit of what a polynomial time verifier can hope to be convinced of?

Ben-Or, Goldwasser, Kilian and Wigderson [BGKW] proposed an extended definition of interactive proof-systems which they called multi-prover interactive proofs. Instead of one prover attempting to convince a verifier that  $x$ , the input

string, is in  $L$ , the prover consists of  $k$  separate agents (or rather  $k$  provers) who jointly attempt to convince a verifier that  $x$  is in  $L$ . The  $k$  provers can cooperate and communicate between them to decide on a common optimal strategy before the interaction with the verifier starts. But, once they start to interact with the verifier, they can no longer send each other messages or see the messages exchanged between the verifier and the “other prover”.

As in single prover interactive proofs, the verifier is a probabilistic polynomial time algorithm, and can exchange upto a polynomial number of messages with each one of the provers (with no restriction on interleaving the exchanged messages) before deciding to accept or reject the input  $x$ .

**Definition.** Let  $P_1, P_2, \dots, P_k$  be Turing machines which are computationally unbounded and  $V$  be a probabilistic polynomial time Turing machine. All machines have a read-only input tape, a work tape and a random tape. All machines share the same input tape. The  $P_i$ 's share an infinite read-only random tape of 0's and 1's. Every  $P_i$  has one write-only communication tape and one write-only communication tape.  $V$  has  $k$  write-only communication tapes and  $k$  read only communication tapes. The  $i$ th write-only communication tape of  $V$  is  $P_i$ 's read-only communication tape, and vice versa. Namely, on its  $i$ th communication tape,  $V$  writes messages to  $P_i$ . We call  $(P_1, P_2, \dots, P_k, V)$  a *k-prover interactive protocol*.

**Definition.** We say that  $L \subset \{0, 1\}^*$  has a *k-prover interactive proof-system* (IPS) if there exists an interactive probabilistic Turing machine  $V$  such that:

- 1  $\exists P_1, P_2, \dots, P_k$  such that  $(P_1, P_2, \dots, P_k, V)$  is a *k-prover interactive protocol* and  $\forall x \in L, \text{prob}(V \text{ accepts input } x) \geq \frac{2}{3}$ .
- 2  $\forall P_1, P_2, \dots, P_k$  such that  $(P_1, P_2, \dots, P_k, V)$  is a *k-prover interactive protocol*,  $\text{prob}(V \text{ accepts input } x) \leq \frac{1}{3}$ .

We let  $\mathbf{IP}_k$  denote the class of languages which are accepted by *k*-prover interactive proof-systems and  $\mathbf{IP}_k[t(n)]$  denote the class of languages which are accepted by *k*-prover interactive proof-system using  $t(n)$  rounds.

It can be shown that as in the single-prover case, for any  $k$ , if  $L \in \mathbf{IP}_k$  then there exists an an interactive proof  $(P_1, P_2, \dots, P_k, V)$  for  $L$  such that for all  $x \in L$ , the  $\Pr(V \text{ accepts } x) = 1$ . The number of rounds, size, and error probability of a multi-prover interactive proof is defined similarly to interactive proofs.

## 10. The Complexity of Multi-Prover Interactive Proofs

Is  $\mathbf{IP}_k$  be different than  $\mathbf{IP}$ ? The hope is that the inability of the provers to communicate with each other will allow the verifier to check the consistency of the provers responses against each other. Intuitively, this additional check makes the verifier trust the answers of the prover more than in the single-prover case. This potentially may allow the verifier to be convinced of membership in harder languages.

It was shown in the original paper of [BGKW] that languages accepted by interactive proofs with a polynomial number of provers are equal to languages accepted by interactive proofs with two provers.

**Theorem 12** [BGKW]. *For all  $k \geq 2$ , if  $L \in IP_k$  then  $L \in IP_2$ .*

Shortly after the results showing  $IP = PSPACE$ , Babai, Fortnow and Lund [BFL] proved that the class of languages having two-prover interactive proof systems is exactly nondeterministic exponential time ( $NEXPTIME$ ).

**Theorem 13** [BFL].  $IP_2 = NEXPTIME$ .

This represents a further step demonstrating the unexpected power of randomization and interaction in efficient provability. It follows that multiple provers with coins are strictly stronger than without, since  $NEXPTIME \neq NP$ . In particular, for the first time, provably polynomial time intractable languages turn out to have “efficient proof systems” since  $NEXPTIME \neq P$ .

The first part of the proof extends the techniques of [LFKN] and [Sh] for the single prover case. The second part is a verification scheme for the multilinearity of an  $n$ -variable function held by the two provers.

### 10.1 Interaction in Multi-Prover Interactive Proofs

Unbounded number of rounds of interaction seemed necessary to realize the full power of interactive proofs. In contrast, work of Fortnow, Rompel and Sipser [FRS], Cai, Condon and Lipton [CCL], and Kilian shows that a two prover proof system can accept every language in  $IP_2$  using bounded number of rounds and achieving a constant error probability.

**Theorem 14** [FRS, CCL], kilian. *For any constant  $\varepsilon$ , any  $L \in IP_2$  has a bounded round two-prover interactive proof with error probability  $\varepsilon$ .*

## 11. Approximating the Clique Problem

A new connection between approximation of combinatorial graph problems and multi-prover interactive proofs was recently discovered.

It is well known that the Clique problem (i.e., finding the size of the largest complete subgraph in a given graph  $G$ ) is  $NP$ -complete. However, the related problem of *approximating* this size (say within a constant factor) is not known to have an efficient algorithm, nor is it known to be  $NP$ -complete. The best known polynomial time approximation algorithm for clique is a factor  $O(\frac{n}{\log^2 n})$  below the optimal [Boppana et al.]

Using the characterization of  $IP_2 = NEXPTIME$ , Feige, Goldwasser, Lovasz, and Safra [FGLS] showed the following.

**Theorem 15** [FGLS]. *If approximating the clique problem within factor  $2^{\log^c n}$  for some  $c < 1$  can be done in time  $2^{\log^k n}$  for some  $k > 1$  then*

1. EXPTIME = NEXPTIME

2.  $NP \subseteq \bigcup_{k>0} DTIME(2^{\log^k n})$

## References

- [AGH] W.Aiello, S. Goldwasser, J. Hastad: On the power of interaction. Proc. 27th IEEE Foundations of Computer Science Conf. 1986 pp. 368–379
- [BFL] L. Babai, L. Fortnow, C. Lund: Non-deterministic exponential time has two-prover interactive protocols. Proceedings, 31st Annual IEEE Symp. on Foundations of Computer Science, 1990, pp. 16–25.
- [Ba] L. Babai: Trading group theory for randomness. Proc. 17th ACM Symp. on Theory of Computing 1985 pp. 421–429
- [Ba2] L. Babai: E-mail and the unexpected power of interaction. University of Chicago Tech Report CS90-15, April 24, 1990
- [BHZ] R. Boppana, J. Hastad, S. Zachos: Does co-NP have short interactive proofs? Information Processing Letters **25** (1987) 127–132
- [BM] L. Babai, S. Moran: Arthur-Merlin Games: a randomized proof system, and a hierarchy of complexity classes. To appear in J. Computer Science and Systems.
- [C] S. Cook, The Complexity of Theorem-Proving Procedures. Proc. of 3rd Symposium of Theory of Computation, 1971 pp. 151–158
- [BGG] M. Ben-Or, O. Goldreich, S. Goldwasser, J. Håstad, J. Kilian, S. Micali, P. Rogaway: IP is in zero-knowledge. Proceedings, Advances in Cryptology, Crypto 1988
- [BGKW] M. Ben-Or, S. Goldwasser, J. Kilian, A. Wigderson: Multi-prover interactive proof systems: Removing intractability assumptions. Proceedings of the 20th STOC, ACM, 1988
- [CCL] J. Cai, A. Condon, R. Lipton: On bounded round multi-prover interactive proof systems. STACS90
- [FRS] L. Fortnow, J. Rompel, M. Sipser: On the power of multi-prover interactive protocols. Proceedings, Structures **88**, pp. 156–161
- [FGLS] L. Lovasz, S. Goldwasser, U. Feige, S. Safra: On the difficulty approximating the clique function. In preparation
- [GMR] S. Goldwasser, S. Micali, and C. Rackoff: The knowledge complexity of interactive proof-systems. Proceedings of the 17th STOC, ACM, 1985, pp. 291–304
- [GMS] O. Goldreich, Y. Mansour, and M. Sipser: Interactive proof systems: Provers that never fail and random selection. Proc. FOCS87
- [GMW] Goldreich, Oded, Silvio Micali, Avi Wigderson: Proofs that yield nothing but the validity of the assertion, and a methodology of cryptographic protocol design. Proceedings of the 27th FOCS, IEEE, 1986, pp. 174–187
- [GS] S. Goldwasser and M. Sipser: On public versus private coins in interactive proof systems. Proc. 18th Symp. on Theory of Computing (1986) pp. 59–68
- [IY] R. Impagliazzo, M. Yung: Direct Minimum Knowledge Computations. Proceedings, Advances in Cryptology, Crypto 1987
- [LFKN] C. Lund, L. Fortnow, H. Karloff, N. Nisan: Algebraic methods for interactive proof systems. Proceedings, 31st Annual IEEE Symp. on Foundations of Computer Science, 1990, pp. 2–10
- [L] C. Lautemann: BPP and the polynomial hierarchy. Info. Proc. Letters **17**, no. 4 (1983) 215–217
- [Sh] A. Shamir: IP=PSPACE. Proceedings, 31st Annual IEEE Symp. on Foundations of Computer Science, 1990, pp. 11–15



# Information Theoretic Reasons for Computational Difficulty

Avi Wigderson \*

The Hebrew University, Jerusalem, Israel and  
Princeton University, Princeton, NJ 08544, USA

## 1. Introduction

We give an intuitive account of the concepts in the title, by considering the following simple number-theoretic example. Imagine two distant players who communicate by exchanging binary messages (bits). One player is given a *prime* number  $x$ , and the second a *composite* number  $y$ , where  $x, y < 2^n$ . The players' task is to find a prime number  $p$ , with  $p < 2n$ , such that  $x \not\equiv y \pmod{p}$ . The existence of such a small prime  $p$  is guaranteed by the prime number theorem and the Chinese remainder theorem.

The players agree beforehand on a “protocol” for exchanging messages. The protocol dictates to each player what message to send at each point, based on his input and the messages he received so far. It also dictates when to stop, and how to determine the answer from the information received. There is no limit on the computational complexity of these decisions, which are free of charge. The cost of the protocol is the number of bits they have to exchange on the *worst case* choice of inputs. We shall be interested in the cost of the best protocol under this measure, which we denote by  $t(n)$ .

There is a trivial protocol in which one player sends his input to the second ( $n$  bits), who computes the answer and sends it ( $\log n$  bits) back to the first. This shows that  $t(n) \leq n + \log n$ .

How small can  $t(n)$  be? Is it possible that <sup>1</sup>  $t(n) = O(\log n)$ , which is (essentially) the trivial lower bound? At present, these trivial upper and lower bounds are the best known!

Why should anyone take the time to think about this problem, besides its innocently simple statement and the challenge of the exponential gap in our knowledge? The reason is that this *information theoretic* problem encodes the *computational difficulty* of primality testing! Answering it is extremely important for computational number theory and theoretical computer science as follows:

---

\* This research was partially supported by the American-Israeli Binational Science Foundation grant number 87-00082.

<sup>1</sup> For two functions  $f(n), g(n)$  we write  $f(n) = O(g(n))$  (resp.  $f(n) = \Omega(g(n))$ ) if there is a fixed positive constant that bounds the ratio  $f(n)/g(n)$  from above (resp. from below) for every integer  $n$ .

If  $t(n) = O(\log n)$ , then for every  $n$  there is a polynomial (in  $n$ ) size logical formula of  $n$  boolean variables, that evaluates to ‘true’ if and only if the values of these variables represent (in binary) a prime integer! This means that primality can be tested highly efficiently, both sequentially and in parallel, which has far reaching practical consequences (e.g. in cryptography and coding).

If  $t(n) \neq O(\log n)$ , the primality function has no such formulæ, and it would be the first explicit example of such a function. This would resolve a close cousin of the  $P$  vs.  $NP$  question, which is perhaps the single most important problem of theoretical computer science. (In a nutshell, this question asks if there exist mathematical theorems whose proofs are much harder to find than to verify.) Essentially no progress was made on this problem since it was first posed informally in a letter of Gödel to von Neumann in the 50s, and formulated by Edmonds, Cook and Levin in the 70s.

There is nothing special about choosing the primality function. If, instead, the number  $x$  given to the first player is a perfect square, and  $y$  is a non-square, then the communication complexity  $t(n)$  of the same communication problem would correspond to the formula complexity of testing if an integer (given in binary) is a perfect square. For this problem it turns out that  $t(n) = O(\log n)$ .

It is not hard to see, by a simple counting argument, that for almost all Boolean functions on the integers, the communication complexity is as bad as can be, i.e.  $t(n) = \Omega(n)$  (again, for input integers less than  $2^n$ ), and the associated Boolean formulae computing them must be of exponential size. But no explicit example of a function requiring superpolynomial formulae size is known.

The computational complexity of Boolean functions, in the model of Boolean circuits and formulæ, called *Circuit Complexity*, was initiated by Shanon and others in the 40s. Its importance stems from the direct connection between resources in this model, circuit size and circuit depth ( $\approx$  the logarithm of formula size), and the “standard” resources time and space (respectively) in Turing machines. Non-trivial results exist mainly for restricted circuits, especially monotone (no negations allowed) circuits. Comprehensive texts on circuit complexity are the books of Wegener [W] and Dunne [D]. [BS] is a beautiful and concise account of the most important issues.

The information theoretic model of *Communication Complexity* was introduced by Yao [Y] in ’79. In his model the two players have to compute a Boolean function of their two inputs (in contrast to the “search” problem defined above). This model was extensively studied in the last ten years, and many significant questions about it were resolved. No complete survey of this area exists, but many of the issues and results are discussed and referenced in [BFS].

In 1988 Karchmer and Wigderson [KW] suggested to study the communication complexity of relations (search problems), and showed how to associate with an arbitrary Boolean function a relation (in essentially the way described above), such that the circuit depth of the given function is exactly the communication complexity of the associated relation. They also showed how to capture in a similar fashion monotone computation. This connection established that

computational difficulty (at least of the “depth” resource), stems from purely information theoretic limitations.

The elegance of the communication model, and the existence of non-trivial machinery to handle it, enabled the sequence of recent lower bounds on the monotone circuit depth of several important functions [KW, RW2, RW1]. We have believe that this approach will be instrumental in breaking the ice and finally obtaining a non-trivial lower bound for some explicit function on general (non-monotone) circuits!

In Section 2 we give the neccesary background from circuit complexity and communication complexity. In Section 3 we describe the connection between the two, and in Section 4 we state the lower bounds obtained via this connection. In Section 5 we give some concluding remarks.

## 2. Preliminaries

### 2.1 Boolean Functions and Circuit Complexity

For an integer  $r$  let  $[r]$  denote the set  $\{1, 2, \dots, r\}$ . The set  $\{0, 1\}^n$  will denote all binary sequences of length  $n$ . We interpret such a sequence as truth assignment to  $n$  Boolean variables  $x_1, x_2, \dots, x_n$ . These variables will encode some object, e.g. a number, a graph, a matrix in a natural way. Any function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is a **Boolean function**. We will be interested in asymptotic complexity of functions, so  $n$  should be thought of as a parameter describing “input size”, and the function  $f$  as a sequence of functions,  $\{f_n\}$ , one for each value of  $n$ . Here are some functions we will deal with in this paper. Let  $x = (x_1, x_2, \dots, x_n) \in \{0, 1\}^n$ .

- $\text{parity}(x) = 1$  iff  $\sum_{i=1}^n x_i$  is odd.
- $\text{majority}(x) = 1$  iff  $\sum_{i=1}^n x_i > n/2$ .
- $\text{prime}(x) = 1$  iff  $x$  is the binary representation of a prime integer.
- $\text{stconn}$ , the st-connectivity function. Here  $x$  represents a graph (binary relation) on a set  $V$ , with  $|V| = m$ ,  $n = m^2$ . So  $x$  can be thought of as an  $m * m$  matrix  $A$  with  $A_{ij} = 1$  ( $i, j \in V$ ) iff there is a direct connection (an *edge*) from  $i$  to  $j$ . For fixed  $s, t \in V$ ,  $\text{stconn}(x) = 1$  iff there is a path from  $s$  to  $t$  in this graph (i.e. if  $s$  and  $t$  are related in the reflexive transitive closure of the relation  $A$ ).
- $\text{clique}$ . Here  $x$  represents an undirected graph on  $m$  vertices, so  $n = m^2$  and  $x_{ij} = x_{ji}$ . A clique is a subset  $U$  of the vertices such that every pair of distinct vertices  $i, j \in U$  is connected by an edge. Now  $\text{clique}(x) = 1$  iff the graph represented by  $x$  has a clique of size (say)  $m/10$ .
- $\text{matching}$ , the perfect matching function. Again  $x$  represents an  $m * m$  matrix  $A$  with rows  $B$  and columns  $G$  with  $A_{bg} = 1$  if boy  $b \in B$  and girl  $g \in G$  are willing to be matched. Now  $\text{matching}(x) = 1$  iff all boys and girls can be simultaneously matched. Equivalently,  $\text{matching}(x) = 1$  iff the permanent of the matrix  $A$  is positive.

There is a natural partial order on the set of sequences  $\{0, 1\}^n$ , namely for sequences  $x, y$  we say that  $x \leq y$  if for all  $i \in [n]$   $x_i \leq y_i$ . A Boolean function

$f$  is called **monotone** if  $f(x) \leq f(y)$  whenever  $x \leq y$ . For example, the functions *majority*, *stconn*, *clique* and *matching* above are monotone, but *parity* and *prime* are not.

A **Boolean circuit** over  $\{x_1, x_2, \dots, x_n\}$  is a sequence of functions  $g_1, g_2, \dots, g_s$  such that for every  $k \in [s]$  either:

1.  $g_k = x_i$  for some  $i \in [n]$ .
2.  $g_k = \bar{x}_i$  (the negation of  $x_i$ ) for some  $i \in [n]$ .
3.  $g_k = g_i \vee g_j$  for some  $i, j < k$ .
4.  $g_k = g_i \wedge g_j$  for some  $i, j < k$ .

Functions of type 1 and 2 are called inputs of the circuit, those of type 3 and 4 are called gates, and the last function  $g_s$  is called the output of the circuit. Thus, a circuit computes a Boolean function of the variables in a natural way — its output. The **size** of the circuit is  $s$ , the number of functions. The **depth** is the length of the longest “path” from an input to the output. More precisely, the depth is defined to be  $d(g_s)$ , where  $d(g_k) = 0$  if  $g_k$  is an input, and  $d(g_k) = 1 + \max(d(g_i), d(g_j))$  if  $g_k$  is a gate.

The **size** of a function  $f$ , denoted  $s(f)$ , is the size of the smallest circuit computing  $f$ . The **depth** of  $f$ , denoted  $d(f)$ , is the depth of the shallowest circuit computing  $f$ . In a vague sense which can be formalized, the size corresponds to the sequential time to compute  $f$ , and the depth to the “parallel” time, as well as space, needed to compute  $f$ . The relationship between these two fundamental complexity measures is far from understood. It is trivial that  $\log s(f) \leq d(f) \leq s(f)$  for every function  $f$ . A nontrivial construction of Paterson and Valiant [PV] slightly improves the right inequality to  $d(f) \leq s(f)/\log s(f)$ .

If we disallow negated inputs, i.e. functions of type 2 in our circuits, it is clear that they can now compute only monotone functions, and indeed we call such circuits **monotone circuits**. Every monotone function can be computed by a monotone circuit. We let  $s^m(f)$  denote the smallest size of a monotone circuit for  $f$ , and similarly with  $d^m(f)$  for monotone depth. The same weak relationship between monotone size and depth as above is the best that is known. To summarize:

**Theorem 1** [PV].

$$\log s(f) \leq d(f) \leq s(f)/\log s(f) \quad (1)$$

$$\log s^m(f) \leq d^m(f) \leq s^m(f)/\log s^m(f). \quad (2)$$

Finally, a **Boolean formula** is a circuit in which every function  $g_i$  can appear at most once on the right hand side of type 3 or 4. In other words, in a formula subfunctions need to be recomputed, and thus it looks like a tree. The **formula size** of a function  $f$ , denoted  $L(f)$ , is the size of the smallest formula for  $f$ , and similarly  $L^m(f)$  stands for the monotone formula size. The depth of a circuit and a formula is clearly the same. Formula size and (circuit) depth are closely related:

**Theorem 2** [Sp].

$$d(f) = \Theta(\log L(f)) \quad (3)$$

$$d^m(f) = \Theta(\log L^m(f)). \quad (4)$$

## 2.2 Communication Complexity

Let  $X, Y, Z$  be finite sets, and let  $R \subseteq X \times Y \times Z$  be a relation. A *deterministic communication protocol*  $A$  over  $(X, Y, Z)$  specifies the exchange of information bits by two players,  $I$  and  $II$ , that initially receive as inputs  $x \in X$  and  $y \in Y$ , respectively. The protocol dictates who sends the first message, and how a player computes his next message from his input and the previously received messages. It also dictates when they terminate, and how they compute a value  $A(x, y) \in Z$ . There is no limit on the complexity of these computations.

Denote by  $c_A(x, y)$  the number of bits exchanged by  $I$  and  $II$  on the input pair  $(x, y)$  when using protocol  $A$ . Let  $c_A(R) = \max_{(x,y) \in X \times Y} c_A(x, y)$ .

We say that  $A$  computes  $R$  if for all  $(x, y) \in X \times Y$  we have  $(x, y, A(x, y)) \in R$ . Then the *deterministic communication complexity* of the relation  $R$  is  $c(R) = \min\{c_A(R) | A \text{ computes } R\}$ . When  $R$  is clear from the context, we sometimes use  $c(X', Y')$  for the communication complexity of  $R$  restricted to the subdomain  $X' \times Y'$ , where  $X' \subseteq X$  and  $Y' \subseteq Y$ .

The original model for communication complexity, as introduced by Yao [Y], dealt only with the computation of functions by the two players (i.e. for every pair of inputs there is a unique output). We accommodate this common notion of computing functions (rather than relations) in a natural way. If  $R$  has the property that for every  $(x, y)$  there is a unique  $z(x, y) \in Z$  with  $(x, y, z) \in R$ , then we identify the relation  $R$  with the function  $R : X \times Y \rightarrow Z$  where  $R(x, y) = z(x, y)$ .

The main advantage in generalising the model to compute relations, is that now we can relate it to the circuit model.

## 3. Boolean Relations and Circuit Depth

### 3.1 Unrestricted Computation

The study of relations (search problems) was initiated by [KW], who showed that the circuit depth of a boolean function is exactly captured by the communication complexity of a related relation.

Specifically, for  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  define the relation

$$R_f \subset f^{-1}(1) \times f^{-1}(0) \times [n]$$

by  $(x, y, i) \in R_f$  iff  $x_i \neq y_i$ . In words, player  $I$  gets an input  $x \in \{0, 1\}^n$  with  $f(x) = 1$ , player  $II$  gets  $y \in \{0, 1\}^n$  with  $f(y) = 0$ , and their task is to find a coordinate where their inputs differ. Note that this search problem is somewhat different than the one defined in the introduction, where the players' task was to find a small prime rather than a coordinate. However, this is just a matter of representation of integers (either binary notation or modular notation), and since there is an efficient conversion between the two, they are essentially equivalent. The heart of this research is the following theorem.

**Theorem 3** [KW]. *For every function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  we have*

$$d(f) = c(R_f). \tag{5}$$

The simple proof of this theorem follows from the fact that there is a one-to-one correspondance between formulae for  $f$  and communication protocols for the relation  $R_f$ . We illustrate this by an example. Consider the parity function *parity* on  $n$  variables. We will assume that  $n$  is a power of two, and denote the complement of this function by  $\overline{\text{parity}}$ . Also, for  $x \in \{0, 1\}^n$ , let  $x_L, x_R \in \{0, 1\}^{n/2}$  denote respectively the left and right half-portions of the sequence  $x$ .

A simple recursive construction of a formula for  $\text{parity}(x)$  (and  $\overline{\text{parity}}(x)$ ) by de-Morgan's rules) is as follows:

If  $n > 1$  then  $\text{parity}(x) = (\text{parity}(x_L) \wedge \overline{\text{parity}}(x_R)) \vee (\overline{\text{parity}}(x_L) \wedge \text{parity}(x_R))$ .

If  $n = 1$  then  $x$  is a single variable  $x_i$ , and  $\text{parity}(x_i) = x_i$ .

Reducing problem size by a factor of two costs depth two ( $\vee$  of  $\wedge$ 's) so this formula has depth  $2 \log n$  (and size  $n^2$ ).

A protocol for  $R_{\text{parity}}$  can also be described recursively. Recall that by definition player *I* has an odd sequence  $x$  ( $\text{parity}(x) = 1$ ), and player *II* has an even sequence  $y$ , and they want to find a coordinate where  $x$  and  $y$  differ.

If  $n > 1$  then player *I* sends the bit  $\text{parity}(x_L)$  and player *II* responds by sending to player *I* the bit  $\text{parity}(y_L)$ . If these two bits are different, then the players can discard the right portion, and recursively find a coordinate in which  $x_L$  and  $y_L$  differ. If the two bits are the same, this implies  $\text{parity}(x_R) \neq \text{parity}(y_R)$ , and they can recursively continue with the right portions.

If  $n = 1$ , then the index of this bit is the required coordinate.

Again, with two bits the problem size is reduced by a factor of two, so the communication complexity is  $2 \log n$ .

Both the formula and the protocol can be viewed as binary trees. In the formula the boolean values propagate from the inputs to the output via the Boolean gates at the nodes. In the protocol the players move along the tree from output to leaves according to their input values. These two object are essentially the same, if we syntactically identify player *I* with tree nodes labeled by an  $\vee$  gate, and player *II* with the  $\wedge$  gate. The formal proof of Theorem 3 below (in two lemmas) follows this simple idea.

**Lemma 4.** *For all functions  $f$  and all  $B_0, B_1 \subseteq \{0, 1\}^n$  such that  $B_0 \subseteq f^{-1}(0)$  and  $B_1 \subseteq f^{-1}(1)$  we have*

$$C(B_1, B_0) \leq d(f). \quad (6)$$

*Proof.* By induction on  $d(f)$ .

If  $d(f) = 0$  then  $f$  is either  $x_i$  or  $\bar{x}_i$ . In either case, we have that for all  $x \in B_1$  and  $y \in B_0$ ,  $x_i \neq y_i$  so that  $i$  is always an answer and  $C(B_1, B_0) = 0$ .

For the induction step we suppose that  $f = f_1 \wedge f_2$  (the case  $f = f_1 \vee f_2$  is treated similarly) so that  $d(f) = \max(d(f_1), d(f_2)) + 1$ . Let  $B_0^j = B_0 \cap f_j^{-1}(0)$  for  $j = 1, 2$ . By induction we have that  $C(B_1, B_0^j) \leq d(f_j)$  for  $j = 1, 2$ . Consider the following protocol for  $B_1$  and  $B_0$ : *II* sends a 0 if  $y \in B_0^1$ , otherwise he sends a 1; the players then follow the best protocol for each of the subcases. We have

$$C(B_1, B_0) \leq 1 + \max_{j=1,2}(C(B_1, B_0^j)) \leq 1 + \max_{j=1,2}(d(f_j)) = d(f). \quad (7)$$

The converse is as follows, and its similar proof is omitted.

**Lemma 5.** *Let  $B_0, B_1 \subseteq B_n$  such that  $B_0 \cap B_1 = \emptyset$ . Then, there exists a function  $f$  with  $B_0 \subseteq f(0)$  and  $B_1 \subseteq f(1)$  such that*

$$d(f) \leq C(B_1, B_0). \quad (8)$$

### 3.2 Monotone Computation

We conclude this section with the analog of Theorem 3 for monotone computation. Let  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  be a monotone function, and define the relation

$$R_f^m \subset f^{-1}(1) \times f^{-1}(0) \times [n]$$

by  $(x, y, i) \in R_f$  iff  $x_i > y_i$ . In words, player  $I$  gets an input  $x \in \{0, 1\}^n$  with  $f(x) = 1$ , player  $II$  gets  $y \in \{0, 1\}^n$  with  $f(y) = 0$ , and their task is to find a coordinate  $i$  where  $x_i = 1$  and  $y_i = 0$ . The existence of such a coordinate is guaranteed by the monotonicity of  $f$ . This relation captures exactly the depth of monotone circuits for  $f$ .

**Theorem 6** [KW]. *For every monotone function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  we have*

$$d^m(f) = c(R_f^m). \quad (9)$$

The proof is essentially the same as the previous one, with a difference in the base case of the induction — the leaves of the tree — which is the only place monotone and nonmonotone computations differ.

It is very convenient to reformulate Theorem 6 in the following way. A **min-term** (or **minimal 1-witness**) of a monotone function  $f$  is a minimal subset  $S \subseteq [n]$  such that any input  $x$  with  $x_i = 1$  for all  $i \in S$  satisfies  $f(x) = 1$ . Let  $\text{MIN}(f)$  be the set of all minterms of  $f$ . It is easy to see that  $\text{MIN}(f)$  characterizes  $f$ .

Similarly, one can define a **maxterm** (or **minimal 0-witness**) as any minimal subset of the variables which, if set to 0, force the function value to 0. Let  $\text{MAX}(f)$  denote the set of all maxterms of  $f$ .

It is easy to see that every minterm must intersect every maxterm (otherwise  $f$  can be simultaneously be forced to be 0 and 1 on the same input). Also, monotone computation of  $f$  does not become easier if we restrict ourselves to inputs which are defined only by minterms and maxterms. Thus the task of the players in  $R_f^m$  turns out to be: Player  $I$  gets a minterm  $S \in \text{MIN}(f)$ , player  $II$  gets a maxterm  $T \in \text{MAX}(f)$ , and they must find an element of the intersection  $S \cap T$ . And Theorem 6 becomes:

**Theorem 7.** *For every monotone function  $f$  we have*

$$d^m(f) = c(\text{MIN}(f), \text{MAX}(f)). \quad (10)$$

As an example, consider the *majority* function on  $n$  bits, say with  $n = 2k - 1$ , an odd number. Then

$$\text{MIN}(\text{majority}) = \text{MAX}(\text{majority}) = \{S \in [n] : |S| = k\}.$$

Clearly, by the pigeonhole principle every minterm and maxterm intersect, and the players should find a member of the intersection. A recursive procedure similar to the one described above yields only an upper bound of  $O((\log n)^2)$  bits (the reader may wish to find this protocol). It is a remarkable result, which was given two different beautiful proofs, that one can do much better!

**Theorem 8** [AKS, V].

$$d^m(\text{majority}) = O(\log n).$$

Even though the same upper bound holds for the communication complexity of the above relation, the two proofs were given in terms of circuits. They are, in some inherent way, construct the circuit in a bottom-up fashion (from inputs to output), and it is hard to translate them into a simple protocol. In the remaining sections we will see that in other cases it is much easier to handle the communication problem, working top-down on the protocol, especially for lower bound purposes.

## 4. Applications

### 4.1 Size vs. Depth

Shanon, who initiated circuit complexity over 40 years ago, observed that simple counting arguments show that most functions are difficult to compute.

**Theorem 9** [Sh]. *Almost all Boolean functions on  $n$  variables  $f$  satisfy*

$$s(f) = \Omega(2^n/n) \tag{11}$$

$$d(f) = \Omega(n). \tag{12}$$

On the other hand, there is no known sequence of functions (that can be explicitly described)  $f$  for which  $s(f) \neq O(n)$  or  $d(f) \neq O(\log n)$ . Finding such a function is a major problem of theoretical computer science.

With this problem seeming too difficult for present techniques, the monotone analogue of this problem was attacked. Nevertheless, explicit monotone functions  $f$  for which  $s^m(f)$  is super linear or  $d^m(f)$  is super logarithmic were not discovered till 1985. That year, in a breakthrough paper Razborov [R1] proved a **super-polynomial** monotone size lower bound on the clique function above, which is known to be in the class  $NP$ .

**Theorem 10** [R1].

$$s^m(\text{clique}) = n^{\Omega(\log n)}.$$

Soon after, exponential lower bounds were given for functions in  $NP$  by Andreev [A] and by Alon and Boppana [AB]. In particular, the later improved the clique lower bound to

**Theorem 11** [AB].

$$s^m(\text{clique}) = \exp(n^{1/6}).$$

These results, together with theorem 1 immediately imply a super logarithmic lower bound on monotone depth.

**Corollary 12.**

$$d^m(\text{clique}) = \Omega(n^{1/6}).$$

As mentioned in the discussion preceding theorem 1 , it is not known if the logarithmic relationship between size and depth is tight. This means that we are looking for a depth lower bound which is superlogarithmic in the size. The first application of the communication complexity viewpoint enabled to find such a result for the *st-connectivity* function, *stconn*.

There is a natural monotone circuit for this function, that is based on raising the matrix describing the graph to the  $m$ th power, by repeatedly squaring it  $\log m$  times. Each such step takes size  $O(m^3) = O(n^{3/2})$  and depth  $O(\log m) = O(\log n)$ , so we get  $s^m(\text{stconn}) = O(n^{3/2} \log n)$ , and  $d^m(\text{stconn}) = O((\log n)^2)$ . Thus, in this circuit, depth is quadratic in the logarithm of the size. Karchmer and Wigderson [KW] proved that this depth bound is optimal regardless of size.

**Theorem 13** [KW].

$$d^m(\text{stconn}) = \Omega((\log n)^2).$$

The proof of this theorem draws intuition from the communication protocol for  $R_{\text{stconn}}^m$  which corresponds to the above circuit. It uses probabilistic arguments to show that after  $O(\log n)$  bits of communication sent by the players in any protocol, one can interpret the remaining steps of the protocol as solving an *st-connectivity* problem of roughly half the size. For more detail, we refer the reader to the original paper. We just note here, that this direction, working from the output of the circuit down to the inputs is natural in the communication complexity framework, but rare in existing circuit lower bounds.

Together with (4) of Theorem 2, we conclude that there is a superpolynomial gap between circuit and formula size in the monotone model.

**Corollary 14.**

$$L^m(\text{stconn}) = n^{\Omega(\log n)}.$$

We end this section by noting that it is not known if this superpolynomial gap can be improved to exponential (for any monotone function, not necessarily explicitly given).

## 4.2 Monotone vs. Nonmonotone

Once we are able to prove monotone lower bounds, it is interesting to ask if negations can help when computing a monotone function. The *clique* function is unlikely to provide the answer, as being *NP*-complete, we expect an exponential size lower bound even for nonmonotone circuits.

The perfect matching function *match*, on the other hand, has a simple polynomial time algorithm and therefore a polynomial size (nonmonotone) circuit. Razborov [R2] was able to apply his methods to this function as well, and prove the first superpolynomial gap between monotone and nonmonotone computation.

**Theorem 15** [R2].

$$s^m(\text{match}) = n^{\Omega(\log n)}.$$

The methods of Alon and Boppana [AB] which improved the clique lower bound to exponential failed to improve the perfect matching bound. However, E. Tardos [T] found that their arguments would work for another function, which while being very similar to *clique*, has a polynomial time algorithm, and thus established an exponential gap between monotone and nonmonotone size.

The analogous question for depth was resolved last year by Raz and Wigderson [RW2]. Again, the function that exemplifies the gap is the perfect matching function *match*. A sequence of results on this problem and parallel algorithms for it led Borodin, von zur Gathen and Hopcroft [BGH] to the following depth upper bound.

**Theorem 16** [BGH].

$$d(\text{match}) = O((\log n)^2).$$

Note that the monotone depth lower bound that follows from taking the logarithm in Theorem 15 only matches this upper bound. However, the communication complexity approach provided a depth bound which is independent of the size:

**Theorem 17** [RW2].

$$d^m(\text{match}) = \Omega(m) = \Omega(\sqrt{n}).$$

The proof of this theorem utilizes the communication complexity approach in a different way. We show that the communication problem  $R_{\text{match}}^m$  encodes another communication problem, called set disjointness. In this problem, each of the two players is given a subset of  $\{1, 2, \dots, m\}$ , and their task is to find if these two subsets are disjoint or not.

More precisely, any protocol for  $R_{\text{match}}^m$  using  $c$  communication bits would lead to a **probabilistic** protocol for set disjointness which uses only  $O(c)$  bits on average. However, a difficult result of Kalyanasundaram and Schnitger [KS] shows that this is possible only if  $c = \Omega(m)$ , from which our theorem follows.

For the exact definitions of probabilistic communication complexity, as well as the probabilistic *reduction* which establish the relationship above we refer the reader to [RW2].

We note that a straightforward reduction from the *clique* to *match* establishes a better depth lower bound for *clique* than the one which follows from the size lower bound of Alon and Boppana in Corollary 12.

**Corollary 18** [RW2].

$$d^m(\text{clique}) = \Omega(m) = \Omega(\sqrt{n}).$$

## 5. Conclusions

We feel that the communication complexity approach to circuit lower bounds was fruitful for a few reasons. The lack of computation in that model. The ability to view computations top-down; this simplifies several known upper and lower bounds (see [K]). And last, but not least, the existing techniques and results for this model that can be used via reductions.

This paper was not intended as a complete survey, and indeed some new results were recently obtained. However, there are still no nontrivial nonmonotone lower bounds, and we believe that the communication complexity approach will play a role in changing this state of affairs!

*Acknowledgements.* The work described in this paper was done in an exciting and enlightening collaboration with Mauricio Karchmer and Ran Raz.

## References

- [A] Andreev, A.E.: On a method for obtaining lower bounds on the complexity of individual monotone functions. *Dokl. Ak. Nauk. SSSR* **282** (1985) 1033–1037 (in Russian). [English translation in: *Sov. Math. Dokl.* **31** (1985) 530–534]
- [AB] Alon, A., Boppana, R.: The monotone circuit complexity of Boolean functions. *Combinatorica* **7**, no. 1 (1987) 1–22
- [AKS] Ajtai, M., Komlós, J., Szemerédi, E.: An  $O(n \log n)$  sorting network. *Proceedings 15th STOC* 1983, pp. 1–9
- [BFS] Babai, L., Frankel, P., Simon, J.: Complexity classes in communication complexity theory. *Proceedings of the 27th FOCS* 1988, pp. 337–347
- [BGH] Borodin, A., von zur Gathen, J., Hopcroft, J.: Fast parallel matrix and GCD Computations. *Proceedings of the 23rd STOC* 1982, pp. 65–71
- [BS] Boppana, R., Sipser, M.: The complexity of finite functions. MIT/LCS/TM-405 (1989)
- [D] Dunne, P.E.: The complexity of Boolean networks. Academic Press, 1988
- [K] Karchmer, M.: Communication complexity — A new approach to circuit depth. MIT Press, 1989
- [KS] Kalyanasundaram, B., Snitger, G.: The probabilistic communication complexity of set intersection. *Proceedings Structure in Complexity Theory* 1987, pp. 41–49
- [KW] Karchmer, M., Wigderson, A.: Monotone circuits for connectivity require super-logarithmic depth. *Proceedings of the 20th STOC* 1988, pp. 539–550
- [PV] Paterson, M., Valiant, L.G.: Circuit size is non-linear in depth. *TCS* **2** (1976) 397–400
- [R1] Razborov, A.A.: Lower bounds for the monotone complexity of some boolean functions. *Dokl. Ak. Nauk. SSSR* **281** (1985) 798–801 (in Russian) [English transl. in: *Sov. Math. Dokl.* **31** (1985) 354–357]

- [R2] Razborov, A.A.: Lower bounds on the monotone network complexity of the logical permanent. *Mat. Zametki* **37** (1985) 887–900 (in Russian) [English transl. in: *Math. Notes of the Academy of Sciences of USSR* **37** (1985) 485–493]
- [RW1] Raz, R., Wigderson, A.: Probabilistic communication complexity of boolean relations. *Proc. of the 30th FOCS* 1989
- [RW2] Raz, R., Wigderson, A.: Monotone circuits for matching require linear depth. *Proceedings of the 22th STOC* 1990
- [Sh] Shannon, C.: The synthesis of two-terminal switching circuits. *Bell Syst. Techn. J.* **28** (1949) 59–98
- [Sp] Spira, P.M.: On time-hardware complexity tradeoffs for Boolean functions. *Proceedings of 4th Hawaii Symp. on System Sciences* 1971, pp. 525–527
- [T] Tardos, E.: The gap between monotone and non-monotone circuit complexity is exponential. *Combinatorica* **8** (1988) 141–142
- [V] Valiant, L.G.: Short monotone formulæ for the majority function. *J. Algorithms* **5** (1984) 363–366
- [W] Wegner, I.: *The complexity of boolean functions*. John Wiley, 1988.
- [Y] Yao, A.C.-C.: Some complexity questions related to distributive computing. *Proceedings of 11th STOC* 1979, pp. 209–213

# Recent Developments in Shock-Capturing Schemes

*Ami Harten*

School of Mathematical Sciences, Tel Aviv University, Ramat Aviv, Tel Aviv, Israel, and  
Department of Mathematics, University of California, Los Angeles, CA 90024-1555, USA

**Abstract.** In this paper we review the development of the shock-capturing methodology, paying special attention to the increasing nonlinearity in its design and its relation to interpolation. It is well-known that high-order approximations to a discontinuous function generate spurious oscillations near the discontinuity (Gibbs phenomenon). Unlike standard finite-difference methods which use a fixed stencil, modern shock-capturing schemes use an adaptive stencil which is selected according to the local smoothness of the solution. Near discontinuities this technique automatically switches to one-sided approximations, thus avoiding the use of discontinuous data which brings about spurious oscillations.

## 1. Introduction

In this paper, we describe and analyze numerical techniques that are designed to approximate weak solutions of hyperbolic systems of conservation laws in several space dimensions. For sake of exposition, we shall describe these methods as they apply to the pure initial value problem (IVP) for a one-dimensional scalar conservation law

$$u_t + f(u)_x = 0, \quad u(x, 0) = u_0(x). \quad (1.1)$$

To further simplify our presentation, we assume that the flux  $f(u)$  is a convex function, i.e.,  $f''(u) > 0$  and that the initial data  $u_0(x)$  are piecewise smooth functions which are either periodic or of compact support. Under these assumptions, no matter how smooth  $u_0$  is, the solution  $u(x, t)$  of the IVP (1.1) becomes discontinuous at some finite time  $t = t_c$ . In order to extend the solution for  $t > t_c$ , we introduce the notion of weak solutions, which satisfy

$$\frac{d}{dt} \int_a^b u \, dx + f(u(b, t)) - f(u(a, t)) = 0 \quad (1.2a)$$

for all  $b \geq a$  and  $t \geq 0$ . Relation (1.2a) implies that  $u(x, t)$  satisfies the PDE in (1.1) wherever it is smooth, and the Rankine-Hugoniot jump relation

$$f(u(y+0, t)) - f(u(y-0, t)) = [u(y+0, t) - u(y-0, t)] \frac{dy}{dt} \quad (1.2b)$$

across curves  $x = y(t)$  of discontinuity.

It is well-known that weak solutions are not uniquely determined by their initial data. To overcome this difficulty, we consider the IVP (1.1) to be the vanishing viscosity limit  $\varepsilon \downarrow 0$  of the parabolic problem

$$(u^\varepsilon)_t + f(u^\varepsilon)_x = \varepsilon u_{xx}^\varepsilon \quad u^\varepsilon(x, 0) = u_0(x), \quad (1.3a)$$

and identify the unique “physically relevant” weak solution of (1.1) by

$$u = \lim_{\varepsilon \downarrow 0} u^\varepsilon. \quad (1.3b)$$

The limit solution (1.3) can be characterized by an inequality that the values  $u_L = u(y - 0, t)$ ,  $u_R = u(y + 0, t)$  and  $s = dy/dt$  have to satisfy; this inequality is called an entropy condition; admissible discontinuities are called shocks. When  $f(u)$  is convex, this inequality is equivalent to Lax’s shock condition

$$a(u_L) > s > a(u_R) \quad (1.4)$$

where  $a(u) = f'(u)$  is the characteristic speed (see [8] for more details).

We turn now to describe finite difference approximations for the numerical solution of the IVP (1.1). Let  $v_j^n$  denote the numerical approximation to  $u(x_j, t_n)$  where  $x_j = jh$ ,  $t_n = nt$ ; let  $v_h(x, t)$  be a globally defined numerical approximation associated with the discrete values  $\{v_j^n\}$ ,  $-\infty < j < \infty, n \geq 0$ .

The classical approach to the design of numerical methods for partial differential equations is to obtain a solvable set of equations for  $\{v_j^n\}$  by replacing derivatives in the PDE by appropriate discrete approximations. Therefore, there is a conceptual difficulty in applying classical methods to compute solutions which may become discontinuous. Lax and Wendroff [9] overcame this difficulty by considering numerical approximations to the *weak formulation* (1.2a) rather than to the PDE (1.1). For this purpose, they have introduced the notion of schemes in conservation form:

$$v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+\frac{1}{2}} - \bar{f}_{j-\frac{1}{2}}) \equiv (E_h \cdot v^n)_j; \quad (1.5a)$$

here  $\lambda = \tau/h$  and  $\bar{f}_{i+\frac{1}{2}}$  denotes

$$\bar{f}_{i+\frac{1}{2}} = f(v_{i-k+1}^n, \dots, v_{i+k}^n); \quad (1.5b)$$

$\bar{f}(w_1, \dots, w_{2k})$  is a numerical flux function which is consistent with the flux  $f(u)$ , in the sense that

$$\bar{f}(u, u, \dots, u) = f(u); \quad (1.5c)$$

$E_h$  denotes the numerical solution operator. Lax and Wendroff proved that if the numerical approximation converges boundedly almost everywhere to some function  $u$ , then  $u$  is a weak solution of (1.1), i.e., it satisfies the weak formulation (1.2a). Consequently discontinuities in the limit solution automatically satisfy the Rankine-Hugoniot relation (1.2b). We refer to this methodology as shock-capturing (a phrase coined by H. Lomax).

In the following, we list the numerical flux function of various 3-point schemes ( $k = 1$  in (1.5b)):

(i) The Lax-Friedrichs scheme [7]

$$\bar{f}(w_1, w_2) = \frac{1}{2}[f(w_1) + f(w_2) - \frac{1}{\lambda}(w_2 - w_1)] \quad (1.6)$$

(ii) Godunov's scheme [1]

$$\bar{f}(w_1, w_2) = f(V(0; w_1, w_2)); \quad (1.7a)$$

here  $V(x/t; w_1, w_2)$  denotes the self-similar solution of the IVP (1.1) with the initial data

$$u_0(x) = \begin{cases} w_1 & x < 0 \\ w_2 & x > 0 \end{cases}. \quad (1.7b)$$

(iii) The Cole-Murman scheme [12]:

$$\bar{f}(w_1, w_2) = \frac{1}{2}[f(w_1) + f(w_2) - |\bar{a}(w_1, w_2)|(w_2 - w_1)] \quad (1.8a)$$

where

$$\bar{a}(w_1, w_2) = \begin{cases} \frac{f(w_2) - f(w_1)}{w_2 - w_1} & \text{if } w_1 \neq w_2 \\ a(w_1) & \text{if } w_1 = w_2 \end{cases}. \quad (1.8b)$$

(iv) The Lax-Wendroff scheme [9]:

$$\bar{f}(w_1, w_2) = \frac{1}{2}\{f(w_1) + f(w_2) - \lambda a\left(\frac{w_1 + w_2}{2}\right)[f(w_2) - f(w_1)]\}. \quad (1.9)$$

Let  $E(t)$  denote the evolution operator of the exact solution of (1.1) and let  $E_h$  denote the numerical solution operator defined by the RHS of (1.5a). We say that the numerical scheme is  $r$ -th order accurate (in a pointwise sense) if its local truncation error satisfies

$$E(\tau) \cdot u - E_h \cdot u = O(h^{r+1}) \quad (1.10)$$

for all sufficiently smooth  $u$ ; here  $\tau = O(h)$ . If  $r > 0$ , we say that the scheme is consistent.

The schemes of Lax-Friedrichs (1.6), Godunov (1.7), and Cole-Murman (1.8) are first order accurate; the scheme of Lax-Wendroff (1.9) is second order accurate.

We remark that the Lax-Wendroff theorem states that if the scheme is convergent, then the limit solution satisfies the weak formulation (1.2b); however, it need not be the entropy solution of the problem (see [4]). It is easy to see that the schemes of Cole-Murman (1.8) and Lax-Wendroff (1.9) admit a stationary “expansion shock” (i.e.,  $f(u_L) = f(u_R)$  with  $a(u_L) < a(u_R)$ ) as a steady solution. This problem can be easily rectified by adding sufficient numerical dissipation to the scheme (see [11] and [3]).

## 2. Interpolatory Schemes and Linear Discontinuities

Let us consider the constant coefficient case  $f(u) = au, a = \text{const.}$  in (1.1), i.e.,

$$u_t + au_x = 0, \quad u(x, 0) = u_0(x), \quad (2.1a)$$

the solution to which is

$$u(x, t) = u_0(x - at). \quad (2.1b)$$

In this case the schemes (1.6) – (1.9) take the form

$$v_j^{n+1} = \sum_{\ell=-K^-}^{K^+} C_\ell(v) v_{j+\ell}^n \equiv (E_h \cdot v^n)_j \quad (2.2)$$

where  $v = \lambda a$  is the CFL number. The coefficients  $C_\ell(v)$  are independent of the numerical solution  $v^n$ ; this makes  $E_h$  a linear operator.

We say that the numerical scheme  $E_h$  is (linearly) stable if

$$\|(E_h)^n\| \leq C \quad \text{for } 0 \leq n\tau \leq T, \quad \tau = O(h). \quad (2.3a)$$

In the constant coefficient case the scheme is stable if and only if it satisfies von Neumann's condition

$$\left| \sum_{\ell=-K^-}^{K^+} C_\ell(v) \ell^{i\xi} \right| \leq 1 \quad \text{for all } 0 \leq \xi \leq \pi. \quad (2.3b)$$

It is easy to see that all the 3 point schemes (1.6) – (1.9) are stable under the CFL condition

$$|v| = |\lambda a| \leq 1. \quad (2.3c)$$

The notion of stability (2.3a) is related to convergence through Lax's equivalence theorem, which states that a consistent linear scheme is convergent if and only if it is stable (see [13] for more details).

Let us denote by  $S_i^r$  the stencil of  $(r+1)$  successive points starting with  $x_i$

$$S_i^r = \{x_i, x_{i+1}, \dots, x_{i+r}\}, \quad (2.4a)$$

let  $P(x; S_i^r, u)$  denote the unique polynomial of degree  $r$  interpolating the  $(r+1)$  values of  $u$  on this stencil and let  $Q(x; u)$  denote the piecewise polynomial interpolation of  $u$

$$Q(x; u) = P(x; S_{i(j)}^r, u) \quad x_{j-1} \leq x \leq x_j. \quad (2.4b)$$

We refer to the numerical scheme

$$v_j^{n+1} = Q(x_j - at; v^n) \quad (2.4c)$$

as interpolatory scheme. Clearly, the interpolatory scheme (2.4) is  $r$ -th order accurate. When  $Q(x; v^n)$  is the piecewise linear interpolation of  $v^n$  (i.e.,  $r = 1, i(j) = j - 1$  in (2.4b)) then (2.4c) is the first-order accurate upwind scheme; in

the constant coefficient case this scheme is identical to those of Godunov (1.7) and Cole-Murman (1.8).

Next let us assume  $a > 0$  and consider the second order case  $r = 2$  in which  $Q(x; v^n)$  is a piecewise-parabolic interpolation of  $v^n$ . There are two different choices of stencil in (2.4): Taking  $Q$  in  $[x_{j-1}, x_j]$  to be the parabola through  $S_{j-1}^2 = \{x_{j-1}, x_j, x_{j+1}\}$  (i.e.,  $i(j) = j - 1$ ) results in the Lax-Wendroff scheme (1.9); taking  $Q$  in  $[x_{j-1}, x_j]$  to be the parabola through  $S_{j-2}^2 = \{x_{j-2}, x_{j-1}, x_j\}$  (i.e.,  $i(j) = j - 2$ ) results in the second-order upwind scheme.

We turn now to consider the application of these schemes to the step function  $H(x)$

$$H(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}, H_j = \begin{cases} 0 & j \leq 0 \\ 1 & j \geq 1 \end{cases}. \quad (2.5a)$$

For the first order upwind scheme we get that

$$Q(x; H) = \begin{cases} 0 & x \leq 0 \\ x/h & 0 \leq x \leq h \\ 1 & h \leq x \end{cases}; \quad (2.5b)$$

for the Lax-Wendroff scheme

$$Q(x; H) = \begin{cases} 0 & x \leq -h \\ \frac{1}{2} \frac{x}{h} (1 + \frac{x}{h}) & -h \leq x \leq 0 \\ 1 - \frac{1}{2} (1 - \frac{x}{h}) (2 - \frac{x}{h}) & 0 \leq x \leq h \\ 1 & h \leq x \end{cases}; \quad (2.5c)$$

for the second order upwind scheme we get that

$$Q(x; H) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{2} \frac{x}{h} (1 + \frac{x}{h}) & 0 \leq x \leq h \\ 1 + \frac{1}{2} (\frac{x}{h} - 1) (2 - \frac{x}{h}) & h \leq x \leq 2h \\ 1 & 2h \leq x \end{cases}. \quad (2.5d)$$

We observe that  $Q$  in (2.5b) is a monotone function of  $x$ ; consequently the numerical solution by Godunov's scheme to these data is also monotone. On the other hand  $Q$  for the second order schemes (2.5c) – (2.5d) is not a monotone function. For the Lax-Wendroff scheme  $Q$  is negative in  $-h \leq x \leq 0$  and has a minimum of  $-0.125$ ; similarly for the second order upwind scheme  $Q$  is larger than 1 in  $h \leq x \leq 2h$  with a maximum of  $1.125$ . This observation explains the Gibbs-like phenomenon of generating spurious oscillations in calculating discontinuous data with these second order schemes.

We say that the scheme  $E_h$  is monotonicity preserving if

$$v \text{ monotone} \Rightarrow E_h \cdot v \text{ monotone}. \quad (2.6)$$

Clearly the numerical solution of a monotonicity preserving scheme to initial data of a step-function is always monotone and therefore the discontinuity propagates without generating spurious oscillations.

Godunov has shown that the *linear* scheme (2.2) is monotonicity preserving if and only if

$$C_\ell(v) \geq 0, \quad -K_- \leq \ell \leq K_+; \quad (2.7)$$

this implies that a monotonicity-preserving scheme which is linear is necessarily only first-order accurate. It took some time to realize the Godunov's monotonicity theorem does not mean that there are no high-order accurate monotonicity preserving schemes; it only means that there are no such linear ones. Hence high-order accurate monotonicity-preserving schemes are nonlinear in an essential way.

The second-order accurate schemes mentioned above are linear because the choice of the stencil (2.4) is fixed. Let us consider now a piecewise-quadratic interpolation which is made nonlinear by an adaptive selection of the stencil in (2.4b). For the interval  $[x_{j-1}, x_j]$  let us consider the two stencils  $S_{j-2}^2 = \{x_{j-2}, x_{j-1}, x_j\}$  and  $S_{j-1}^2 = \{x_{j-1}, x_j, x_{j+1}\}$ , and select the one in which the interpolant is smoother. If we measure the smoothness of  $u$  by the second derivative of the corresponding parabola we select

$$i(j) = \begin{cases} j-2 & \text{if } |\frac{d^2}{dx^2} P(x; S_{j-2}^2, u)| \leq |\frac{d^2}{dx^2} P(x; S_{j-1}^2, u)| \\ j-1 & \text{otherwise} \end{cases}. \quad (2.8a)$$

When we apply this selection of stencil to the step-function  $H(x)$  (2.5a) we get that for  $[x_{-1}, x_0]$  we choose the stencil  $S_{-2}^2 = \{x_{-2}, x_{-1}, x_0\}$  for which  $P(x; S_{-2}^2, H) \equiv 0$ ; for the interval  $[x_1, x_2]$  we choose the stencil  $S_1^2 = \{x_1, x_2, x_3\}$  for which  $P(x; S_1^2, H) \equiv 1$ . As is evident from comparing (2.5c) and (2.5d) it does not matter which stencil we assign to  $[x_0, x_1]$  since both parabolae are monotone there; with (2.8a) we select  $S_{-1}^2$  for  $[x_0, x_1]$ . Thus we get in (2.4)

$$Q(x; H) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{2} \frac{x}{h} (1 + \frac{x}{h}) & 0 \leq x \leq h \\ 1 & h \leq x \end{cases} \quad (2.8b)$$

which is a monotone function of  $x$  although it is actually a piecewise-quadratic polynomial.

The use of an adaptive stencil is the main idea behind the Essentially Non-Oscillatory (ENO) schemes to be described later in this paper. It extends to high order of accuracy in a straightforward manner: For  $r$ -th order accuracy we consider for  $[x_{j-1}, x_j]$  the  $r$  stencils  $S_{j-r}^r, S_{j-r+1}^r, \dots, S_{j-1}^r$ . We choose  $i(j)$  in (2.4b) to be the one which minimizes

$$\left| \frac{d^r}{dx^r} P(x; S_i^r, u) \right| \quad \text{for } i = j-r, \dots, j-1. \quad (2.9)$$

### 3. Total Variation Stability and TVD Schemes

An immense body of work has been done to find out whether stability of constant coefficient scheme with respect to all “frozen coefficients” associated with the problem, implies convergence in the variable coefficient case and in the nonlinear case.

In the variable coefficient case, where the numerical solution operator is linear and Lax’s equivalence theorem holds, it comes out that the stability of the variable coefficient scheme depends strongly on the dissipativity of the constant coefficient one, i.e., on the particular way it damps the high-frequency components in the Fourier representation of the numerical solution.

In the nonlinear case, under assumptions of sufficient smoothness of the PDE, its solution and the functional definition of the numerical scheme, Strang proved that linear stability of the first variation of the scheme implies its convergence; we refer the reader to [13] for more details.

In the case of discontinuous solutions of nonlinear problems, linearly stable schemes are not necessarily convergent; when such a scheme fails to converge, we refer to this case as “nonlinear instability.” The occurrence of a nonlinear instability is usually associated with insufficient numerical dissipation which triggers exponential growth of the high-frequency components of the numerical solution.

The following theorem states that a stronger sense of stability, namely uniform boundedness of the total variation of the numerical solution, does imply convergence to a weak solution.

**Theorem 3.1.** *Let  $v_h$  be a numerical solution of a conservative scheme (1.5).*

(i) *If*

$$\text{TV}(v_h(\cdot, t)) \leq C \cdot \text{TV}(u_0) \quad (3.1)$$

*where  $\text{TV}(\cdot)$  denotes the total variation in  $x$  and  $C$  is a constant independent of  $h$  for  $0 \leq t \leq T$ , then any refinement sequence  $h \rightarrow 0$  with  $\tau = O(h)$  has a convergent subsequence  $h_j \rightarrow 0$  that converges in  $L^{\text{loc}}$  to a weak solution of (1.1).*

(ii) *If  $v_h$  is consistent with an entropy inequality which implies uniqueness of the IVP (1.1), then the scheme is convergent (i.e., all subsequences have the same limit, which is the unique entropy solution of the IVP (1.1)).*

We say that the scheme  $E_h$  is Total Variation Diminishing (TVD) if

$$\text{TV}(E_h \cdot v) \leq \text{TV}(v) \quad (3.2)$$

where

$$\text{TV}(w) = \sum_j |w_{j+1} - w_j|. \quad (3.3)$$

Clearly TVD schemes satisfy (3.1) with  $C = 1$  and therefore are TV stable.

In [2] we have shown that if the scheme can be written in the form

$$v_j^{n+1} = v_j^n + C_{j+\frac{1}{2}}^+ A_{j+\frac{1}{2}} v^n - C_{j-\frac{1}{2}}^- A_{j-\frac{1}{2}} v^n \quad (3.4a)$$

where  $C_{j+\frac{1}{2}}^\pm$  satisfy for all  $j$

$$C_{j+\frac{1}{2}}^\pm \geq 0, \quad C_{j+\frac{1}{2}}^+ + C_{j+\frac{1}{2}}^- \leq 1 \quad (3.4b)$$

then the scheme is TVD; here  $A_{i+\frac{1}{2}}v^n = v_{i+1}^n - v_i^n$ . Applying this lemma to the general scheme

$$v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+\frac{1}{2}} - \bar{f}_{j-\frac{1}{2}}) \quad (3.5a)$$

$$\bar{f}_{j+\frac{1}{2}} = \frac{1}{2}(f_j + f_{j+1} - q_{j+\frac{1}{2}}A_{j+\frac{1}{2}}v^n) \quad (3.5b)$$

we get that if  $\lambda q$  satisfies

$$\lambda|\bar{a}_{j+\frac{1}{2}}| \leq \lambda q_{j+\frac{1}{2}} \leq 1 \quad (3.6a)$$

then the scheme (3.5) is TVD; here

$$\bar{a}_{j+\frac{1}{2}} = \frac{f_{j+1} - f_j}{A_{j+\frac{1}{2}}v}. \quad (3.6b)$$

This shows that the Cole-Murman scheme (1.8) for which  $q = |\bar{a}|$  is TVD subject to the CFL restriction  $\lambda|\bar{a}_{j+\frac{1}{2}}| \leq 1$ .

Using conditions (3.4b) it is possible to construct TVD schemes which are second-order accurate in the  $L_1$ -sense (see [2] and [14]). However, TVD schemes are at most second-order accurate (see [5]). In order to design higher-order accurate shock capturing schemes we introduce the notion of Essentially Non-Oscillatory (ENO) schemes.

## 4. ENO Schemes

In this section we describe high-order accurate Godunov-type schemes which are a generalization of Godunov's scheme (1.7) and van Leer's MUSCL scheme [10].

We start with some notations: Let  $\{I_j\}$  be a partition of the real line; let  $A(I)$  denote the interval-averaging (or "cell-averaging") operator

$$A(I) \cdot w = \frac{1}{|I|} \int_I w(y) dy; \quad (4.1)$$

let  $\bar{w}_j = A(I_j) \cdot w$  and denote  $\bar{w} = \{\bar{w}_j\}$ . We denote the approximate reconstruction of  $w(x)$  from its given cell-averages  $\{\bar{w}_j\}$  by  $R(x; \bar{w})$ . To be precise,  $R(x; \bar{w})$  is a piecewise-polynomial function of degree  $(r-1)$ , which satisfies

$$(i) \quad R(x; \bar{w}) = w(x) + O(h^r) \quad \text{wherever } w \text{ is smooth} \quad (4.2a)$$

$$(ii) \quad A(I_j) \cdot R(\cdot; \bar{w}) = \bar{w}_j \quad (\text{conservation}). \quad (4.2b)$$

Finally, we define Godunov-type schemes by

$$v_j^{n+1} = A(I_j) \cdot E(\tau) \cdot R(\cdot; v^n) \equiv (\bar{E}_h \cdot v^n)_j \quad (4.3a)$$

$$v_j^0 = A(I_j)u_0; \quad (4.3b)$$

here  $E(t)$  is the evolution operator of (1.1).

In the scalar case, both the cell-averaging operator  $A(I_j)$  and the solution operator  $E(\tau)$  are order-preserving, and consequently also total-variation diminishing (TVD); hence

$$\text{TV}(\bar{E}_h \cdot \bar{w}) \leq \text{TV}(R(\cdot; \bar{w})). \quad (4.4)$$

This shows that the total variation of the numerical solution of Godunov-type schemes is dominated by that of the reconstruction step.

We turn now to describe the recently developed essentially non-oscillatory (ENO) schemes of [5, 6], which can be made accurate to any finite order  $r$ . These are Godunov-type schemes (4.3) in which the reconstruction  $R(x; \bar{w})$ , in addition to relations (4.2), also satisfies

$$\text{TV}(R(\cdot; \bar{w})) \leq \text{TV}(\bar{w}) + O(h^{1+p}), \quad p > 0 \quad (4.5)$$

for any piecewise-smooth function  $w(x)$ . Such a reconstruction is essentially non-oscillatory in the sense that it may not have a Gibbs-like phenomenon at jump-discontinuities of  $w(x)$ , which involves the generation of  $O(1)$  spurious oscillations (that are proportional to the size of the jump); it can, however, have small spurious oscillations which are produced in the smooth part of  $w(x)$ , and are usually of the size  $O(h')$  of the reconstruction error (4.2a).

When we use an essentially non-oscillatory reconstruction in a Godunov-type scheme, it follows from (4.4) and (4.5) that the resulting scheme (4.3) is likewise essentially non-oscillatory (ENO) in the sense that for all piecewise-smooth function  $w(x)$

$$\text{TV}(\bar{E}_h \cdot \bar{w}) \leq \text{TV}(\bar{w}) + O(h^{1+p}), \quad p > 0; \quad (4.6)$$

i.e., it is “almost TVD.” Property (4.6) makes it reasonable to believe that the total variation of the numerical solution is uniformly bounded. We recall that by Theorem 3.1, this would imply that the scheme is convergent (at least in the sense of having convergent subsequences). This hope is supported by a very large number of numerical experiments.

Next we describe one of the techniques to obtain an ENO reconstruction. To simplify our presentation we assume that  $\{I_j\}$  is a uniform partition

$$I_j = (x_{j-1}, x_j), \quad x_j = jh.$$

Given cell averages  $\{\bar{w}_j\}$  of piecewise-smooth function  $w(x)$ , we observe that

$$h\bar{w}_j = \int_{x_{j-1}}^{x_j} w(y) dy = W(x_j) - W(x_{j-1}) \quad (4.7a)$$

where

$$W(x) = \int_{x_0}^x w(y) dy \quad (4.7b)$$

is the primitive function of  $w(x)$ . Hence we can easily compute the point values  $\{W(x_j)\}$  by summation

$$W(x_i) = h \sum_{j=i_0}^i \bar{w}_j. \quad (4.7c)$$

Once we have computed the point values of the primitive function we use the ENO interpolation technique (2.4), (2.9) to obtain  $Q(x; W)$ , an  $r$ -th order piecewise-polynomial interpolation of  $W$ , i.e.,

$$Q(x; W) = P(x; S_i^r, W) \quad \text{for } x_{j-1} \leq x \leq x_j \quad (4.8a)$$

where  $P(x; S_i^r, W)$  is the unique  $r$ -th degree polynomial which interpolates  $W$  over the stencil  $S_i^r = \{x_i, x_{i+1}, \dots, x_{i+r}\}$ , and  $i(j)$  is chosen so that

$$\left| \frac{d^r}{dx^r} P(x; S_i^r, W) \right| = \min_{j-r \leq i \leq j-1} \left| \frac{d^r}{dx^r} P(x; S_i^r, W) \right|. \quad (4.8b)$$

We define  $R(x; \bar{w})$  by

$$R(x; \bar{w}) = \frac{d}{dx} Q(x; W). \quad (4.9)$$

We observe that if  $w(x)$  is smooth in  $(x_{j-1}, x_j)$  then for  $h$  sufficiently small the algorithm (4.8b) will select a stencil  $S_i^r$  in which  $w(x)$  is smooth. It follows then from standard interpolation theorems that

$$R(x; \bar{w}) = \frac{d}{dx} P(x; S_i^r, W) = \frac{d}{dx} W + O(h^r) = w(x) + O(h^r) \quad (4.10)$$

which is property (4.2a). Furthermore (4.10) holds in every interval except for those in which  $w(x)$  has a discontinuity. As we have seen in the examples (2.5) and (2.8b) the Gibbs-phenomenon is associated with intervals near the discontinuity and not with the interval that contains the discontinuity. This is why the reconstruction (4.8) – (4.9) satisfies the ENO property (4.5); in [2] we show that the second-order accurate ENO scheme is actually TVD. The conservation property (4.2b) follows directly from the definition (4.9):

$$\begin{aligned} A(I_j)R(\cdot; \bar{w}) &= \frac{1}{h} \int_{x_{j-1}}^{x_j} \frac{d}{dx} Q(x; W) dx = \frac{1}{h} [Q(x_j; W) - Q(x_{j-1}; W)] \\ &= \frac{1}{h} [W(x_j) - W(x_{j-1})] = \bar{w}_j. \end{aligned} \quad (4.11)$$

The abstract scheme (4.3) can be written in the standard conservation form (1.5). To do so let us denote by  $\tilde{v}(x, t)$  the solution in the small of the IVP

$$\begin{cases} (\tilde{v})_t + f(\tilde{v})_x = 0 \\ \tilde{v}(x, 0) = R(x; v^n), \quad 0 \leq t \leq \tau \end{cases} \quad (4.12)$$

and integrate this PDE over  $I_j \times [0, \tau]$ ; using the divergence theorem and (4.2b) we get that  $v^{n+1}$  in (4.3) can be expressed by

$$v_j^{n+1} = v_j^n - \lambda [\hat{f}_{j+\frac{1}{2}} - \hat{f}_{j-\frac{1}{2}}] \quad (4.13a)$$

where

$$\hat{f}_{j+\frac{1}{2}} = \frac{1}{\tau} \int_0^\tau f(\tilde{v}(x_j, t)) dt \quad (4.13b)$$

In the first-order case the scheme (4.13) is identical to Godunov's scheme and the numerical flux (4.13b) can be expressed in a closed form by (1.7b). For higher order schemes we use a numerical flux which is an appropriate approximation to (4.13b) (see [6] for more details).

We remark that the ENO schemes are related to the interpolatory schemes of Sect. 2 as follows: In the constant coefficient case a fixed choice of stencil (i.e.,  $i(j)-j = \text{constant}$  in (4.8a)) results in the interpolatory scheme (2.4) corresponding to the same choice of stencil.

## References

1. Godunov, S. K.: A difference scheme for numerical computation of discontinuous solutions of equations of fluid dynamics. *Math. Sbornik* **47**, 271–306 (1959) (in Russian)
2. Harten, A.: High resolution schemes for hyperbolic conservation laws. *J. Comp. Phys.* **49**, 357–393 (1983)
3. Harten, A., Hyman, J. M.: A self-adjusting grid for the computation of weak solutions of hyperbolic conservation laws. *J. Comp. Phys.* **50**, 235–269 (1983)
4. Harten, A., Hyman, J. M., Lax, P. D.: On finite-difference approximations and entropy conditions for shocks. *Comm. Pure Appl. Math.* **29**, 297–322 (1976)
5. Harten, A., Osher, S.: Uniformly high-order accurate non-oscillatory schemes, I. *SIAM J. Numer. Anal.* **24**, 279 (1987)
6. Harten, A., Engquist, B., Osher, S., Chakravarthy, S. R.: Uniformly high-order accurate non-oscillatory schemes, III. *J. Comp. Phys.* **71**, 231 (1987)
7. Lax, P. D.: Weak solutions of nonlinear hyperbolic equations and their numerical computation. *Comm. Pure Appl. Math.* **7**, 159–193 (1954)
8. Lax, P. D.: Hyperbolic systems of conservation laws and the mathematical theory of shock waves. Society for Industrial and Applied Mathematics, Philadelphia, 1972
9. Lax, P. D., Wendroff, B.: Systems of conservation laws. *Comm. Pure Appl. Math.* **13**, 217–237 (1960)
10. van Leer, B.: Towards the ultimate conservative difference schemes V. A second order sequel to Godunov's method. *J. Comp. Phys.* **32**, 101–136 (1979)
11. Majda, A., Osher, S.: Numerical viscosity and entropy condition. *Comm. Pure Appl. Math.* **32**, 797–838 (1979)
12. Murman, E. M.: Analysis of embedded shock waves calculated by relaxation methods. *AIAA J.* **12**, 626–633 (1974)
13. Richtmyer, R. D., Morton, K. W.: Difference methods for initial value problems. 2nd edn. Interscience-Wiley, New York 1967
14. Sweby, P. K.: High resolution schemes using flux limiters for hyperbolic conservation laws. *SIAM J. Numer. Anal.* **21**, 995–1011 (1984)



# Undirected Multiflow Problems and Related Topics – Some Recent Developments and Results

Alexander V. Karzanov

Institute for System Studies, Academy of Sciences of USSR, 9, Prospect 60 Let Oktyabrya  
117312 Moscow, USSR

**Abstract.** A multiflow (multicommodity flow), arising as a natural extension of the well-known notion of a network flow, is a popular object studied in linear programming and combinatorial optimization. We discuss some recent results on undirected multiflow problems and ideas behind them, concerning combinatorial and computational aspects, such as: (i) the existence of feasible and optimal solutions with small denominators, (ii) special solvability criteria and minimax relations, (iii) efficient solution algorithms, (iv) a relationship between multiflows and packings of cuts and metrics, and some others.

## 1. Preliminaries

**Fractionality.** We start with some basic notion. Suppose that  $\mathcal{K}$  is a collection (a class) of linear programs  $P$  of the form: (i) find  $x \in \mathbf{Q}^n$  satisfying  $Ax \leq b$  (the feasibility problem), or (ii) maximize (or minimize)  $c^T x$  subject to  $Ax \leq b$ ,  $x \in \mathbf{Q}^n$  (the optimization problem), where  $A$  is an integral  $m \times n$ -matrix and  $b$  ( $c$ ) is an integral  $m$ -vector ( $n$ -vector). Define the fractionality  $\varphi(P)$  of  $P$  to be the minimum positive integer  $k$  for which  $P$  has a rational feasible (optimal) solution  $x$  such that  $kx$  is integral; if  $P$  has no feasible (optimal) solution, we put  $\varphi(P) := 1$ . The fractionality  $\varphi(\mathcal{K})$  of  $\mathcal{K}$  is defined as  $\sup\{\varphi(P) | P \in \mathcal{K}\}$ ; if  $\varphi(\mathcal{K}) = \infty$  we say that  $\mathcal{K}$  has unbounded fractionality.

In combinatorial applications, a class  $\mathcal{K}$  usually describes the linear programming relaxation of a certain combinatorial problem, and determining  $\varphi(\mathcal{K})$  can be rather difficult. For example, let  $A_G$  be the  $\{0, 1\}$ -matrix whose rows correspond to the edges  $e \in E$  of a graph  $G = (V, E)$  and the columns are the incidence vectors of the perfect matchings  $M \subseteq E$  in  $G$ , that is, each vertex of  $G$  is covered by exactly one edge of  $M$ . Then the well-known conjecture of Berge-Fulkerson [Fu2] (cl. [Se3]) that “any bridgeless cubic graph  $G$  can be covered by six perfect matchings so that each edge is covered twice” is equivalent to the assertion that  $\varphi(\mathcal{K}) = 2$  where  $\mathcal{K}$  is the set of programs  $\max\{1^T x | A_G x \leq 1, x \geq 0\}$  for such  $G$ s (actually it is even unknown whether  $\varphi(\mathcal{K})$  is finite).

**Flows and Multiflows.** We shall deal only with undirected flows and multiflows. By a graph we mean a finite undirected graph without multiple edges and loops; an edge  $\{u, v\}$  will be denoted by  $uv$ .

Consider a network  $N = (G, H, c)$  consisting of a (supply) graph  $G = (V, E)$ , a (demand) graph  $H = (T, U)$  with  $T \subseteq V$  and a nonnegative integer-valued

vector  $c \in \mathbf{Z}_+^E$  (of edge capacities); we suppose that each vertex of  $H$  is covered by an edge. An  $s - t$  chain in  $G$  is a subgraph  $L = (V_L, E_L)$  of  $G$  with  $V_L = \{s = v_0, v_1, \dots, v_m = t\}$  and  $E_L = \{v_i v_{i+1} | i = 1, \dots, m\}$ . Define  $\mathcal{L} := \mathcal{L}(G, H)$  to be  $\bigcup_{st \in U} \mathcal{L}(G, st)$  where  $\mathcal{L}(G, st)$  is the set of  $s - t$  chains in  $G$ . A *multiflow* in  $N$  is a nonnegative rational-valued function  $f$  on  $\mathcal{L}$  satisfying the capacity constraints:

$$\zeta^f(e) := \sum (f(L)|L \in \mathcal{L}, e \in E_L) \leq c(e) \quad \text{for all } e \in E. \quad (1)$$

For  $st \in U$  the restriction  $f_{st}$  of  $f$  to  $\mathcal{L}(G, st)$  is called a *flow* between  $s$  and  $t$  of *value*  $v(f_{st}) := \sum (f(L)|L \in \mathcal{L}(G, st))$ . The *total* value  $v(f)$  of  $f$  is  $\sum_{st \in U} v(f_{st})$ .

Three types of problems on multiflows are well-known.

- (2) *Demand problem*,  $D(G, H, c, d)$ : given a vector  $d \in \mathbf{Z}_+^U$  (of demands), find  $f$  satisfying  $v(f_{st}) = d(st)$  for all  $st \in U$  (or establish that such an  $f$  does not exist).
- (3) *Maximization problem*,  $M(G, H, c)$ : find  $f$  with  $v(f)$  maximum (a *maximum multiflow*).
- (4) *Minimum cost multiflow problem*,  $M(G, H, c, a)$ : given a vector  $a \in \mathbf{Z}_+^E$  (of edge costs), find a *maximum* multiflow  $f$  whose *total cost*  $\sum_{e \in E} a(e) \zeta^f(e)$  is minimum (clearly this is reduced to a linear program).

An obvious necessary condition of solvability of (2) (that is, the existence of a required  $f$ ) is the *cut condition*

$$\Delta(c, d, X) := c(\delta_G(X)) - d(\delta_H(X \cap T)) \geq 0 \quad \text{for all } X \subseteq V, \quad (5)$$

where  $\delta(X') = \delta_G(X')$  denotes the set of edges of a graph  $G' = (V', E')$  with one endvertex in  $X' \subseteq V'$  and the other in  $V' \setminus X'$  (a “cut” in  $G'$ ); and for  $c' \in \mathbf{Q}^{E'}$  and  $E'' \subseteq E'$ ,  $c'(E'')$  stands for  $\sum (c'(e)|e \in E'')$ . When  $|U| = 1$  (2)–(4) are specified to be the classical (single-commodity) demand, maximum and minimum cost maximum flow problems. Well-known results on flows are that if  $|U| = 1$  then: (i) (5) is sufficient for solvability of (2), and (ii) (2) has an integral (feasible) solution whenever it is solvable [Me, FF]. Similarly, if  $|U| = 1$  then (3) and (4) have integral optimal solutions. This can be written as  $\varphi(D(H)) = \varphi(M(H)) = \varphi(C(H)) = 1$  for  $|U| = 1$ , where  $D(H)$  is the set of problems  $D(G, H, c, d)$  with fixed  $H$  and arbitrary  $G$ ,  $c$  and  $d$ ;  $M(H)$  and  $C(H)$  are defined in a similar way.

When  $G$  and  $H$  vary, determining  $\varphi(P)$  for an “individual” problem  $P$  in (2), (3) or (4) seems to be rather difficult. However, as we shall see in Sects. 2–5, this task has been successfully accomplished for many classes  $D(H)$ ,  $M(H)$  and  $C(H)$ . In particular, the values  $\varphi(C(H))$  have been found for all  $H$ , and there is only one  $H$  for which  $\varphi(D(H))$  is still unknown. In Sect. 5 results on the fractionality for planar supply graphs  $G$  are discussed. For some other results on multiflows see [Fr].

The cut condition (5) gives a special solvability criterion for (2) in the case  $|U| = 1$ . This is equivalent to the minimax relation “the maximum value of a flow between vertices  $s$  and  $t$  is equal to the minimum capacity  $c(\delta(X))$  of a cut  $\delta(X)$  separating  $s$  and  $t$ ” ( $\delta(X)$  separates  $s$  and  $t$  if  $|X \cap \{s, t\}| = 1$ ). We shall see that the bounded fractionality behaviour for  $D(H)$  or  $M(H)$  with other  $H$ 's is also accompanied by appearance of certain special solvability criteria or minimax relations.

As to computational aspects, we are not going to present here a survey of various approaches developed to solve multiflow problems. The purpose is only to outline the idea of one general approach, the so-called “splitting-off” method. By applying this method one can prove constructively the existence of an integral or half-integral solutions for certain “difficult” problems of type (2) and (3), and develop efficient algorithms to solve them. Apparently splitting-off techniques appeared originally in [RW] in connection with the two-commodity flow problem.

**Complexity.** We briefly recall the notions of polynomial and strongly polynomial algorithms. For details, see, e.g., [GJ]. An algorithm is said to be *polynomial* if its running time (the number of bit operations) is bounded by a polynomial in the input size of the problem (the *input size* is the number of bits occurring in the input). Speaking of a strongly polynomial algorithm, one means that the input consists of two parts. The first part describes a “combinatorial structure”, and the second part is a list of “numerical data”. The *dimension* of the input is the amount of bits in description of this structure plus the amount of numerical data. For instance, the input size of  $D(G, H, c, d)$  is  $O(|V|^2 + \sum_{e \in E} \log_2(c(e) + 1) + \sum_{u \in V} \log_2(d(u) + 1))$  while its dimension is  $O(|V|^2)$ . A polynomial algorithm is *strongly polynomial* if, roughly speaking, it consists of elementary arithmetic operations and data transfers, the number of these operations is bounded by a polynomial in the dimension of the input, and the size of data arised in the algorithm is polynomially bounded in the input size.

*Remarks.* (i) The above definition of a flow in the “edge-chain” form is known to be equivalent to the usual “node-edge” definition of a flow as a function on directed edges satisfying conservation conditions [FF]. The first form is preferable for us to explain some combinatorial ideas and methods. Though this form leads to appearance of constraint matrices for (2)–(4) with possibly exponential in  $|V|$  number of columns we support explicitly only chains (“columns”) with non-zero values of a multiflow; this provides efficiency of algorithms mentioned below.

(ii) The problems (2)–(4), being stated in the node-edge form, can be solved by general-purpose polynomial algorithms like the ellipsoid method [Kh]. Moreover, there are polynomial algorithms of finding a basis solution [GLS] and even, by a method of Tardos [Ta], strongly polynomial algorithms (since the constraint matrix has entries 0, 1, -1, and hence its size is polynomially bounded in  $|V|$ ). However these general methods do not guarantee (at least for the problems (3) and (4)) that the denominators of the obtained (optimal) solution will not exceed the fractionality of a problem; this will be explained in Sect. 3.

**Related Problems.** There is a special kind of duality between multiflows and packings of cuts or metrics. This can be illustrated by the following example. Let  $l \in \mathbf{Z}_+^E$  be a vector (of edge *lengths*) associated with a connected graph  $G = (V, E)$ . A simple fact is that for  $s, t \in V$  there exist cuts  $\delta(X_1), \dots, \delta(X_k)$  in  $G$  such that: (i)  $|\{i | e \in \delta(X_i)\}| \leq l(e)$  for all  $e \in E$ , and (ii)  $|\{i | \delta(X_i) \text{ separates } s \text{ and } t\}| = \text{dist}_l(s, t)$ . Here  $\text{dist}_l(u, v)$  is the minimum  $l$ -length  $\sum(l(e) | e \in E_L)$  of an  $s - t$  chain  $L$  in  $G$ . The problem of determining cuts satisfying (i)–(ii) is dual, in a sense, to (2) with  $U = \{st\}$ . In Sect. 6 we discuss similar problems on cuts and metrics and their relation to multiflows.

## 2. Demand Problem

Applying Farkas' lemma to (2) and making simple transformations one obtains the following *metrical criterion* [Lo2]:  $D(G, H, c, d)$  is solvable if and only if

$$\sum_{e \in E} c(e)m(e) - \sum_{u \in U} d(u)m(u) \geq 0 \quad (6)$$

holds for each metric  $m$  on  $V$  such that

$$m \text{ is primitive and has an extremal graph } \Gamma \text{ with } E_\Gamma \subseteq U. \quad (7)$$

Here by a *metric* on  $V$  we mean a nonnegative rational-valued function  $m$  on the set of unordered pairs in  $V$  satisfying  $m(xx) = 0$  and  $m(xy) + m(yz) \geq m(xz)$  for any  $x, y, z \in V$  (we use the term "metric" rather than "semimetric");  $m$  is called *primitive* if  $m' + m'' = m$ , where  $m'$  and  $m''$  are metrics, implies  $m' = \lambda m$  for some  $\lambda$ ; an *extremal graph* of  $m$  is a minimal graph  $\Gamma = (V_\Gamma, E_\Gamma)$  with  $V_\Gamma \subseteq V$  such that for any distinct  $u, v \in V$  there is  $st \in E_\Gamma$  for which  $m(st) = m(su) + m(uv) + m(vt)$ .

We say that a metric  $m$  on  $V$  is *induced* by a graph  $Q = (V_Q, E_Q)$  if there is a mapping  $\sigma$  from  $V$  onto  $V_Q$  such that for  $u, v \in V$ ,  $m(uv)$  is equal to the distance in  $Q$  between  $\sigma(u)$  and  $\sigma(v)$ . A simplest example of a primitive metric gives that induced by the graph  $K_2$ , called the *cut metric* ( $K_p$  is the complete graph on  $p$  vertices). Clearly, for a cut metric (6) turns into the inequality in (5).

Papernov found the complete list of demand graphs  $H$  such that for all  $G$  any metric  $m$  as in (7) is a cut metric. By the argument above, this describes the set of  $H$  for which the cut condition (5) gives a criterion of solvability of (2).

**Theorem 1** [Pa]. *Suppose that  $H$  is either  $K_4$  or  $C_5$  or a 2-star. Then (2) is solvable if and only if (5) holds. For any other  $H$  there are  $G$ ,  $c$  and  $d$  such that (5) holds but (2) has no solution.*

(Here  $C_5$  is the circuit on 5 vertices, and a graph is called a *p-star* if there are at most  $p$  vertices in it covering all edges.) Theorem 1 says nothing about numbers  $\varphi(D(H))$  for  $H$  occurring in it. This number gives another theorem, due to Lomonosov, which extends results in [Hu, RW, ADK, Se4]. We say that a vector  $(c, d)$  is *Eulerian* if  $c(\delta_G(X)) + d(\delta_H(X \cap T))$  is even for every  $X \subseteq V$ . Clearly  $\varphi(D(H)) \leq 2\varphi(D^e(H))$  where  $D^e(H)$  is the set of problems in  $D(H)$  with Eulerian  $(c, d)$ .

**Theorem 2.** [Lo1, Lo3]. *If  $H$  is as in Theorem 1,  $(c, d)$  is Eulerian and (5) holds then  $P = D(G, H, c, d)$  has an integral solution, that is,  $\varphi(D^e(H)) = 1$ .*

Theorem 2 implies that  $\varphi(D(H)) \leq 2$  for  $H$  as above (actually  $\varphi(D(H)) = 2$  unless  $H$  is a 1-star). There are strongly polynomial algorithms to find a half-integral or integral (in the Eulerian case) solution when  $H$  is a 2-star [Ch1, Se2] or  $H \in \{K_4, C_5\}$  [Ka1]. These algorithms exploit, in essence, augmentation techniques that construct a solution starting with zero multiflow. We now explain how to derive Theorem 2 directly from Theorem 1 by use of splitting-off operations, as shown in [Ka6, Sc3].

Suppose that (5) holds. Then, by Theorem 1,  $P$  has a solution  $f$ . One may assume that  $d \neq 0$ ,  $c(e) > 0$  for all  $e \in E$ , and  $st \notin E$  whenever  $st \in U$ . For a

pair  $\pi = \{uv, vw\}$  of edges of  $G$  let  $\alpha = \alpha(\pi) \leq \min\{c(uv), c(vw)\}$  be the maximum rational for which  $D(G', H, c', d)$  is still solvable, where  $G'$  and  $c'$  arise from  $G$  and  $c$  by decreasing  $c(uv)$  and  $c(vw)$  and increasing  $c(uw)$  by  $\alpha$  (if  $uw \notin E$  we add the edge  $uw$  of capacity  $\alpha$ ); this operation is called *splitting off*  $\pi$  by  $\alpha$ . Let  $\pi$  be chosen so that  $\alpha$  is maximum. If  $\alpha \geq 1$  we split off  $\pi$  by 1, obtaining a “simpler” solvable problem  $D(G', H, c'', d)$  for which  $(c'', d)$  is again Eulerian, and the result follows by induction. But if  $\alpha < 1$  then, by Theorem 1, there is  $X \subseteq V$  such that  $\Delta'' := \Delta(c'', d, X) < 0$  ( $\Delta(\cdot, \cdot, X)$  is defined in (5)). Since  $\Delta := \Delta(c, d, X)$  and  $\Delta''$  are even,  $\Delta \geq 0$  and, obviously,  $c(\delta(X)) - c''(\delta(X)) \leq 2$ , we have  $\Delta = 0$ , whence  $\alpha = 0$ . This implies  $f = 0$ , and hence  $d = 0$ ; a contradiction.

The above proof can be transformed to a strongly polynomial algorithm as follows. Choose a vertex  $v \in V$  of the current network  $(G, H, c)$ , consider pairs  $\pi = \{uv, vw\}$ , one by one, and split off  $\pi$  by  $\lfloor \alpha(\pi) \rfloor$ . If  $v \in V \setminus T$ , remove  $v$  from  $G$ . Repeat the same for a new vertex  $v'$ , and so on. As a result, one eventually gets  $\tilde{G} = (\tilde{V}, \tilde{E})$  and  $\tilde{c}$  such that  $\tilde{V} = T$  and  $\tilde{c}(e) \geq d(e)$  for  $e \in U$ . Now a required multiflow in the original network is constructed in a natural way by using the obtained numbers  $\alpha$ . One shows that calculation of  $\alpha(\pi)$  can be reduced to solving  $O(1)$  minimum cut problems. This provides a strongly polynomial algorithm.

Now let  $\Gamma_1 + \dots + \Gamma_p$  denote the graph that is the union of disjoint graphs  $\Gamma_1, \dots, \Gamma_p$ . For  $H = K_2 + K_2 + K_2$  and arbitrary  $k \in \mathbb{Z}_+$  one can construct a solvable problem (2) of fractionality at least  $k$  [Lo1, Lo3]. This and simple observation that if  $H'$  is a subgraph of  $H$  then  $\varphi(D(H')) \leq \varphi(D(H))$  imply the following result.

**Theorem 3.** *If  $H$  contains a matching of 3 edges then  $\varphi(D(H)) = \infty$ .*

The only graphs different from those in Theorems 2 and 3 are: (i) certain subgraphs of  $K_5$ , (ii) the union of  $K_3$  and a 1-star, (iii)  $K_3 + K_3$ . The case (ii) is easily reduced to (i). The following theorem generalizes Theorem 2.

**Theorem 4** [Ka6]. *If  $H = K_5$  then  $\varphi(D^e(H)) = 1$ .*

The proof of this theorem given in [Ka6] and based on splitting-off techniques is rather complicated. First, one shows that a metric satisfying (7) for  $H = K_5$  is either a cut metric or a metric induced by  $K_{2,3}$ , called a 2,3-metric ( $K_{p,q}$  is the complete bipartite graph with parts of  $p$  and  $q$  vertices). Thus (2) is solvable if and only if (6) holds for all cut metrics and 2,3-metrics on  $V$ . Second, unlike the proof above (when  $\alpha = \alpha(\pi)$  is always an integer), in our case  $\alpha$  can take half-integer values; in particular,  $\alpha = 1/2$  is possible (in which case the “obstacle”  $m$  violating (6) after splitting off  $\pi$  by 1 is a 2,3-metric). The core of the proof is to show, using combinatorial properties of 2,3-metrics, that if  $\alpha(\pi) < 1$  for all  $\pi$  then  $\alpha$  is 0 everywhere.

This proof also can be turned into a strongly polynomial algorithm (however, using the ellipsoid method). To get such an algorithm, one shows that determining  $\alpha(\pi)$  is reduced to solving  $O(1)$  problems  $P$ : given  $c' \in \mathbb{Z}_+^E$  and a metric  $\varrho$  on  $T$  that is either a cut metric or a 2,3-metric, find a metric  $m$  on  $V$  such that  $m$  coincides with  $\varrho$  on  $T$  and  $\sum(c'(e)m(e)|e \in E)$  is minimum. The size of the constraint matrix for  $P$  is a polynomial in  $|V|$ , hence  $P$  is solvable in a strongly polynomial time, by [Ta].

Theorems 2–4 yield  $\varphi(D(H))$  for all  $H$  except  $H = K_3 + K_3$ . A simple example shows that  $\varphi(D^e(K_3 + K_3))$  is at least 2, and a conjecture is that it is exactly 2.

Theorems 2 and 4 have some consequence in lattice theory. A set  $S = \{a_1, \dots, a_k\} \subseteq \mathbf{Q}^n$  is said to form a *Hilbert basis* if the intersection of the lattice  $a_1\mathbf{Z} + \dots + a_k\mathbf{Z}$  with the cone in  $\mathbf{Q}^n$  generated by  $S$  coincides with  $\{\lambda_1 a_1 + \dots + \lambda_k a_k | \lambda_1, \dots, \lambda_k \in \mathbf{Z}_+\}$  (cf. [GP]). Let  $S(G, H)$  be the set of the following vectors in  $\mathbf{Q}^E \times \mathbf{Q}^U$ : (i)  $(\chi_L, e_{st})$ ,  $L$  is an  $s - t$  chain in  $G$ ,  $st \in U$ ; (ii)  $(\chi_C, 0)$ ,  $C$  is a circuit in  $G$ ; (iii)  $(2e_e, 0)$ ,  $e \in E$ ; here  $\chi_L$  ( $\chi_C$ ) is the incidence vector of  $E_L$  ( $E_C$ ) in  $\mathbf{Q}^E$ , and  $e_e$  ( $e_{st}$ ) is the  $e$ -th ( $st$ -th) unit basis vector in  $\mathbf{Q}^E$  ( $\mathbf{Q}^U$ ). Then  $S(G, H)$  forms a Hilbert basis when  $H$  is a subgraph of  $K_5$ , or a 2-star, or the union of  $K_3$  and a 1-star.

### 3. Maximization Problem

The set of  $H$ 's for which  $M(H)$  has bounded fractionality turns out to be larger than that for  $D(H)$ . The complete list of such  $H$ 's is unknown, but the values  $\varphi(M^*(H))$  have been determined for all  $H$ . Here  $M^*(H)$  is the set of programs  $P = M^*(G, H, c)$  dual to those in  $M(H)$ ;  $P$  can be written as: minimize  $c^T l$  subject to  $l \in \mathbf{Q}_+^E$  and  $\text{dist}_l(s, t) \geq 1$  for each  $st \in U$ .

We say that  $c$  is *inner Eulerian* if  $c(\delta(\{v\}))$  is even for any  $v \in V \setminus T$ . Clearly  $\varphi(M(H)) \leq 2\varphi(M^e(H))$  where  $M^e(H)$  is the set of problems in  $M(H)$  with inner Eulerian  $c$ . The following theorem, due to Lomonosov and the author, generalizes results in [Hu, RW, Lov, Ch2]. Let  $\mathcal{A}(H)$  denote the collection of maximal independent sets (*anticliques*) of  $H$  ( $A \subset T$  is *independent* if  $st \notin U$  for all  $s, t \in A$ ).

**Theorem 5.** [Ka3, Lo3]. *Let  $\mathcal{A}(H)$  have a partition  $\{\mathcal{A}_1, \mathcal{A}_2\}$  such that each  $\mathcal{A}_i$  consists of pairwise disjoint sets (in other words,  $H$  is the complement of the line graph of a bipartite graph). Then  $\varphi(M^e(H)) = 1$ .*

(A weaker statement that  $\varphi(M(H)) \leq 2$  occurred in [KL].) Again, as for Theorems 2 and 4, splitting-off techniques can be applied to prove Theorem 5 [Ka3]. A sketch is as follows. It turns out that the dual problem for  $H$  in question can be solved separately from the primal one. More precisely, one shows that: (i) the dual problem has an optimal solution  $l$  of a special form [KL], namely,

$$l = \frac{1}{2} \sum_{A \in \mathcal{A}(H)} \chi(\delta(X_A)) \quad (8)$$

where  $X_A$ ,  $A \in \mathcal{A}$ , are subsets of  $V$  such that  $X_A \cap T \subseteq A$ , and each  $s \in T$  is contained in just one  $X_A$  ( $\chi(E')$  is the incidence vector of  $E' \subseteq E$ ); and (ii) these  $X_A$ 's can be found by solving the minimum cut problem in a certain network of dimension  $O(|E||V|)$ . One can see that  $c^T l$  is an integer when  $c$  is inner Eulerian. This observation implies that  $\alpha(\pi)$  is a multiple of  $1/2$  for any  $\pi = \{uv, vw\} \subseteq E$ , where  $\alpha(\pi)$  is the maximum of  $\lambda \leq \min\{c(uv), c(vw)\}$  for which splitting off  $\pi$  by  $\lambda$  preserves the maximum value of a multiflow. The central point of the proof is to show that if  $uv, vw, vz$  are three edges with  $\alpha(\{uv, vw\}) = \alpha(\{vw, vz\}) = 1/2$  then  $\alpha(uv, vz) \geq 1$ ; this enables us to apply induction. The proof can be turned into a “pure combinatorial” strongly polynomial algorithm.

The expression (8) shows that  $\varphi(M^*(H)) \leq 2$  for  $H$  as in Theorem 5. The exact values of  $\varphi(M^*(H))$  were pointed out in [Ka7] for all  $H$ . In particular, the following result was stated there. Specify the set of  $H$ 's with the property:

- (9) for any three pairwise intersecting antcliques  $A, B, C$  in  $H$ ,  $A \cap B = B \cap C = C \cap A$ .

**Theorem 6.**  $\varphi(M^*(H)) \in \{1, 2, 4\}$  if  $H$  satisfies (9) and  $\varphi(M^*(H)) = \infty$  otherwise.

The first part can be reformulated in polyhedral terms as follows. Let  $P(G, H)$  be the polyhedron  $\{l \in \mathbf{Q}^E \mid l \geq 0, \text{dist}_i(s, t) \geq 1 \text{ for any } st \in U\}$ . If  $H$  is as in (9) then  $P(G, H)$  is 1/4-integral, that is, each face in it contains a 1/4-integral point.

Another consequence of Theorem 6 is that if  $H$  is not as in (9) then  $\varphi(M(H)) = \infty$  because  $\varphi(M(H')) \geq \varphi(M^*(H'))$  for all  $H'$ . The latter follows from a general statement [Ka7] (extending a result on totally dual integral systems [Fu1, EG]): let  $A$  be a nonnegative  $m \times n$ -matrix,  $b$  be an integral  $m$ -vector, and let the program  $D(c) := \max\{y^T b \mid y \geq 0, y^T A \leq c\}$  have a  $1/k$ -integral optimal solution for every nonnegative integral  $n$ -vector  $c$ ; then the polyhedron  $\{x \in \mathbf{Q}^n \mid x \geq 0, Ax \geq b\}$  is  $1/k$ -integral.

It is unknown whether  $\varphi(M(H))$  is finite for each  $H$  as in (9); a conjecture is that  $\varphi(M(H)) \leq 4$ .

Finally, we show that  $M(G, H, c)$  can have an optimal basis solution  $f$  such that the denominator of some component of  $f$  exceeds  $\varphi(M(H))$ , as it was mentioned in Sect. 1. Take a demand problem  $D = D(G', H', c, d)$  with  $H' = (T', U') = K_2 + K_2 + K_2$  such that  $\varphi(D) > 2$  (existing by Theorem 3). Let  $U = \{s_i t_i \mid i = 1, 2, 3\}$ . Add new vertices  $p_i, q_i$  and edges  $p_i s_i, q_i t_i, i = 1, 2, 3$ , to  $G'$ , forming the graph  $G$ , and put  $c(p_i s_i) := c(q_i t_i) := d(s_i t_i)$ . Let  $H = (T, U)$  be the complete graph on  $T$ , and  $B := \{p_i q_i \mid i = 1, 2, 3\}$ . Then  $\varphi(M(H)) \leq 2$ , by Theorem 5. Now take the objective function  $h(f) := \sum(v(f_u) \mid u \in B)$ . Then any optimal basis solution of  $\max\{h(f) \mid f \text{ a multiflow in } (G, H, c)\}$  is an optimal basis solution of  $M(G, H, c)$  and it determines a solution of  $D(G' H' c, d)$  in a natural way, whence the vector  $2f$  is not integral.

## 4. Minimum Cost Problem

According to the multi-terminal version of the minimum-cost maximum-flow theorem [FF],  $\varphi(C(H)) = 1$  when  $H$  is a complete bipartite graph. The following result was stated in [Ka2].

**Theorem 7.** If  $H$  is a complete  $p$ -partite graph with  $p \geq 3$  (that is,  $\mathcal{A}(H)$  consists of  $p$  pairwise disjoint sets) then  $\varphi(C(H)) = 2$ .

Theorem 7 is a consequence of a pseudo-polynomial algorithm (an algorithm of complexity  $O(c(E)Q(|V|))$ ,  $Q(n)$  is a polynomial in  $n$ ) which finds a half-integral optimal primal solution. This algorithm extends the minimum-cost augmenting path method in [FF] based on ideas of the primal-dual method in linear programming. Recently the author found a strongly polynomial algorithm using a general method in [Ta].

On the other hand, it was shown in [Ka5] that if  $H$  is not a complete  $p$ -partite ( $p \geq 2$ ) then  $\varphi(C(H)) = \infty$ .

## 5. Demand Problem for Planar Graphs

Speaking of a planar graph  $G$ , we mean that  $G$  is explicitly embedded in the plane without intersecting edges. There are known several cases of the demand problem (2) when it is solvable provided that the cut condition (5) holds. The most interesting of them are the following:

(C1) the graph  $(V, E \cup U)$  is planar [Se5];

(C2) there is a set  $\mathcal{H}$  of two faces in  $G$  such that each edge of  $H$  connects vertices in the boundary of some face in  $\mathcal{H}$  [Ok];

(C3)  $U = \{s_1 t_1, \dots, s_k t_k\}$  and there are two inner faces  $I$  and  $J$  in  $G$  so that  $s_1, \dots, s_k$  occur in clockwise order in  $I$  and  $t_1, \dots, t_k$  do so in  $J$  [Sc2].

Moreover, in (C1)–(C3), if  $(c, d)$  is Eulerian and (5) holds then (2) has an integral solution.

Now we consider the case similar to (C2) for  $|\mathcal{H}| \geq 3$ . A simple example with  $G = K_{2,3}$  shows that (5) is, in general, not sufficient for solvability of (2). However, the above result is extended, in a sense, as follows.

**Theorem 8** [Ka10]. *Let  $|\mathcal{H}| = 3$ .*

(i) *(2) is solvable if and only if the metrical condition (6) holds for all  $m$  such that  $m$  is a cut metric or a 2,3-metric on  $V$ .*

(ii) *If  $(c, d)$  is Eulerian and (2) is solvable then (2) has an integral solution.*

It was shown in [Ka10] that if  $|\mathcal{H}| = 4$  (or more) then (ii) is, in general, not true, and there are infinitely many “types” of metrics  $m$  necessary for checking solvability of (2) for all corresponding  $G$  and  $H$ .

The statement (i) follows from a result on packing of cuts and 2,3-metrics (Theorem 10(ii)(b) below). To prove (ii) we use (i) and the splitting-off method as in the proof of Theorem 4. There are certain difficulties when applying this method, because in order to keep planarity we should take only those pairs  $\pi$  of edges of  $G$  which are contained in the boundary of a face of  $G$ . The core of the proof is to show that if  $\alpha(\pi) < 1$  for all such  $\pi$ 's then there are three edges of capacity 1 in  $G$  such that the graph  $G'$  obtained by removing these edges consists of three components, each containing just one face in  $\mathcal{H}$ . Now (ii) is proved by using Okamura's theorem for (C2).

## 6. Packings of Cuts and Metrics

There is a kind of duality that connects solvability conditions for the demand problem (2) with a certain packing problem on metrics. It can be expressed in a general form as follows.

**Proposition 9.** *Given  $G = (V, E)$ ,  $H = (T, U)$  and a set  $M$  of metrics on  $V$ , the following statements are equivalent:*

- (i) *for any  $c$  and  $d$ , (2) is solvable if and only if (6) holds for all  $m \in M$ ;*
- (ii) *for any  $l \in \mathbb{Z}_+^E$ , there exist  $m_1, \dots, m_k \in M$  and  $\lambda_1, \dots, \lambda_k \in \mathbf{Q}_+$  so that:*

$$\lambda_1 m_1(e) + \dots + \lambda_k m_k(e) \leq l(e) \quad \text{for all } e \in E; \quad (10)$$

and

$$\lambda_1 m_1(st) + \dots + \lambda_k m_k(st) = \text{dist}_l(st) \quad \text{for all } st \in U. \quad (11)$$

This is easily proved by applying Farkas' lemma or the cone polarity. Proposition 9 enables us to derive results on packing of metrics directly from corresponding solvability theorems for (2) (like Theorems 1, 4, 8), and vice versa. Note that this relationship gives only theorems on the existence of *rational*  $\lambda$ . There are stronger, integral, version for some of these theorems; as a rule, their proofs are based on special, sometimes complicated, combinatorial approaches. Now we present some results in this area. We say that a vector  $l \in \mathbb{Z}_+^E$  is *bipartite-like* if the  $l$ -length of every circuit in  $G$  is even.

**Theorem 10.** *Let  $l$  be bipartite-like.*

(i) (10) and (11) hold for some cut metrics  $m_i$ 's and integral  $\lambda_i$ 's in the following cases: (a)  $H$  is  $K_4$  or  $C_5$  or a 2-star [Ka4] (cf. [Se1] for  $H = K_2 + K_2$ ); (b)  $G$  and  $H$  are as in (C1) in Sect. 5 [Se5]; (c)  $G$  and  $H$  are as in (C2) in Sect. 5 [Sc1] (see [Ka8] for a strongly polynomial algorithm).

(ii) (10) and (11) hold for  $m_1, \dots, m_k$ , where  $m_i$  is a cut metric or a 2,3-metric, and integers  $\lambda_1, \dots, \lambda_k$  in the following cases: (a)  $H$  is  $K_5$  or the union of  $K_3$  and a 1-star [Ka9]; (b)  $G$  and  $H$  are as in Theorem 8 [Ka10].

There is a connection of the problem (10)–(11) and the problem (P): given a metric  $m$ , decide whether  $m$  is contained in the conic hull of metrics from a certain collection  $M$ . Such a connection was demonstrated in [Ka4] in terms of an extremal graph of  $m$  for  $M$  consisting of the set of cut-metrics. It was also shown there that for this  $M$  the problem (P) (or, equivalent, the problem “whether  $m$  is embeddable isometrically in the space  $L^1$ ” [De]) is NP-hard.

## References

- [ADK] Adelson-Velsky, G.M., Dinitis, E.A., Karzanov, A.V.: Flow Algorithms. Nauka, Moscow, 1975 (Russian)
- [Ch1] Cherkassky, B.V.: A finite algorithm for solving the two-commodity flow problem. Ekon. Mat. Metody **9** (1973) 1147–1149 (Russian)
- [Ch2] Cherkassky, B.V.: A solution of a problem of multicommodity flows in a network. Ekon. Mat. Metody **13** (1) (1977) 143–151 (Russian)
- [De] Deza, M.: On the Hamming geometry of unitary cubes. Dokl. Akad. Nauk SSSR **134** (1960) 1037–1040 (Russian)
- [EG] Edmonds, J., Giles, R.: A min-max relation for submodular function on graphs. In: Hammer, P.L., Johnson, E.L., Korte, B., Nemhauser, G.L. (eds.), Studies in integer programming. (Ann. Discr. Math., vol. 10.) North-Holland, Amsterdam 1977, pp. 185–204
- [Fr] Frank, A.: Packing paths, circuits and cuts, a survey. Report 88532-OR, Inst. für Diskrete Mathematik, Bonn 1988
- [Fu1] Fulkerson, D.R.: Blocking polyhedra. In: Harris, B. (ed.), Graph theory and its applications. Academic Press, New York 1970, pp. 93–112
- [Fu2] Fulkerson, D.R.: Blocking and antiblocking pairs of polyhedra. Math. Progr. **1** (1971) 168–194
- [FF] Ford, L.R., Fulkerson, D.R.: Flows in networks. Princeton Univ. Press, Princeton, NJ, 1962
- [GJ] Garey, M.R., Johnson, D.S.: Computers and intractability: A guide to the theory of NP-completeness. W.H. Freeman, San Francisco, CA, 1979

- [GP] Giles, R., Pulleyblank, W.: Total dual integrality and integral polyhedra. *Linear Algebra and Appl.* **25** (1979) 191–196
- [GLS] Grötschel, M., Lovász, L., Schrijver, A.: Geometric algorithms and combinatorial optimization. (*Algorithms and Combinatorics*, vol. 2.) Springer, Berlin Heidelberg New York, 1988
- [Hu] Hu, T.C.: Multi-commodity network flows. *J. ORSA* **11** (1963) 344–360
- [Ka1] Karzanov, A.V.: Combinatorial methods to solve cut-determined multiflow problems. In: *Combinatorial methods for flow problems*. Inst. for System Studies, Moscow, 1979, iss. 3, pp. 6–69 (Russian)
- [Ka2] Karzanov, A.V.: A minimum cost maximum multiflow problem. In: *Combinatorial methods for flow problems*. Inst. for System Studies, Moscow, 1979, iss. 3, pp. 138–156 (Russian)
- [Ka3] Karzanov, A.V.: On multicommodity flow problems with integer-valued optimal solutions. *Dokl. Akad. Nauk SSSR* **280** (4) (1985) 789–792 (Russian) (English transl.: *Soviet Math. Dokl.* **31** (1) (1985) 151–154)
- [Ka4] Karzanov, A.V.: Metrics and undirected cuts. *Math. Programming* **35** (1985) 183–198
- [Ka5] Karzanov, A.V.: Unbounded fractionality in maximum multiflow and minimum cost multiflow problems. In: Fridman, A.A. (ed.), *Problems of discrete optimization and methods to solve them*. Centr. Econom. and Math. Inst., Moscow, 1987, pp. 123–135 (Russian)
- [Ka6] Karzanov, A.V.: Half-integral five-terminus flows. *Discr. Appl. Math.* **18** (3) (1987) 263–278
- [Ka7] Karzanov, A.V.: Polyhedra related to undirected multicommodity flows. *Linear Algebra and Appl.* **114–115** (1989) 293–328
- [Ka8] Karzanov, A.V.: Packing of cuts realizing distances between certain vertices of a planar graph. *Discr. Math.* **85** (1990) 73–87
- [Ka9] Karzanov, A.V.: Sums of cuts and bipartite metrics. *Europ. J. Comb.* **11** (1990) 473–484
- [Ka10] Karzanov, A.V.: Paths and metrics in a planar graph with three and more distinguished faces, In two parts. Res. Reps. RR 816-M and RR 817-M, IMAG ARTEMIS, Grenoble, 1990 (to appear in *J. Comb. Theory (B)*)
- [KL] Karzanov, A.V., Lomonosov, M.V.: Systems of flows in undirected networks. In: *Mathematical programming etc.* Inst. for System Studies, Moscow, 1978, iss. 1, pp. 59–66 (Russian)
- [Kh] Khachiyan, L.G.: Polynomial algorithms in linear programming. *Zhurnal Vychislitelnoj Matematiki i Matematicheskoi Fiziki* **20** (1980) 53–72 (Russian)
- [Lo1] Lomonosov M.V.: Solutions of two problems on flows in networks. Unpublished manuscript, 1976 (Russian)
- [Lo2] Lomonosov M.V.: On a system of flows in a network. *Problemy Peredatchi Informacii* **14** (1978) 60–73 (Russian)
- [Lo3] Lomonosov M.V.: Combinatorial approaches to multiflow problems. *Discr. Appl. Math.* **11** (1) (1985) 1–94
- [Lov] Lovasz, L.: On some connectivity properties of Eulerian graphs. *Acta Math. Akad. Sci. Hung.* **28** (1976) 129–138
- [Me] Menger, K.: Zur allgemeinen Kurventheorie. *Fundam. Math.* **10** (1927) 96–115
- [Ok] Okamura, H.: Multicommodity flows in graphs. *Discr. Appl. Math.* **6** (1983) 55–62
- [Pa] Papernov, B.A.: On existence of multicommodity flows. In: Fridman, A.A. (ed.), *Studies in discrete optimizations*. Nauka, Moscow, 1976, pp. 230–261 (Russian)
- [RW] Rothschild, B., Whinston, A.: Feasibility of two-commodity network flows. *Oper. Res.* **14** (1966) 1121–1129

- [Sc1] Distances and cuts in planar graphs. Report OS-R8610, Mathematical Centre, Amsterdam, 1986 (to appear in *J. Comb. Th. (B)*)
- [Sc2] Schrijver, A.: The Klein bottle and multicommodity flows. Report OS-R8810, Mathematical Centre, Amsterdam, 1986 (to appear in *Combinatorica*)
- [Sc3] Schrijver, A.: Short proofs on multicommodity flows and cuts. Preprint, 1988 (to appear in *J. Comb. Th. (B)*)
- [Se1] Seymour, P.D.: A two-commodity cut theorem. *Discr. Math.* **23** (1978) 177–181
- [Se2] Seymour, P.D.: A short proof of the two-commodity flow theorem. *J. Comb. Th. (B)* **26** (1979) 370–371
- [Se3] Seymour, P.D.: On multi-colouring of cubic graphs, and conjectures of Fulkerson and Tutte. *Proc. Lond. Math. Soc. (3)* **38** (1979) 423–460
- [Se4] Seymour, P.D.: Four-terminus flows. *Networks* **10** (1980) 79–86
- [Se5] Seymour, P.D.: On odd cuts and planar multicommodity flows. *Proc. Lond. Math. Soc., Ser. III* **42** (1981) 178–192
- [Ta] Tardos, E.: A strongly polynomial algorithm to solve combinatorial linear programs. *Oper. Res.* (1986) 250–256



# Computing Vortex Sheet Motion

Robert Krasny

Department of Mathematics, University of Michigan, Ann Arbor, MI 48109, USA

## 1. Introduction

Coherent vortex structures occur in many types of fluid flow including mixing layers, jets and wakes. A vortex sheet is a mathematical model for such structures, in which the shear layer is approximated by a surface across which the tangential fluid velocity has a jump discontinuity. Vortex sheet motion belongs to the field of vortex dynamics, one of the main approaches to understanding fluid turbulence.

Careful numerical experiments have helped advance the mathematical study of vortex sheets. Difficulties arise in computing vortex sheet motion due to short wavelength instability, singularity formation, and spiral roll-up. This paper reviews the problem of computing vortex sheet motion and presents several applications. See [2] for a sample of other vortex models and numerical methods.

## 2. Analytic Evolution and Singularity Formation

A vortex sheet is defined by a curve  $z(\Gamma, t)$  in the complex plane, where  $\Gamma$  is the circulation parameter and  $t$  is time. The evolution equation is [4,32],

$$\frac{\overline{\partial z}}{\partial t}(\Gamma, t) = \int_{-\infty}^{\infty} K(z(\Gamma, t) - z(\tilde{\Gamma}, t))d\tilde{\Gamma} , \quad K(z) = \frac{1}{2\pi iz} . \quad (1)$$

The Cauchy principal value of the integral is taken. Equation (1) says that a point on the vortex sheet moves with the average of the two limiting velocities, as the curve is approached from either side.

A flat vortex sheet of constant strength  $z(\Gamma, t) = \Gamma$  is an equilibrium solution of (1). Linear stability analysis shows that short wavelength perturbations can grow arbitrarily fast (Kelvin-Helmholtz instability). This means that the linearized initial value problem is ill-posed in the sense of Hadamard. However, Sulem et al. [35] have proven that if the initial perturbation is an analytic function of  $\Gamma$ , then the solution of (1) remains analytic for a positive time interval.

Birkhoff conjectured that instability and nonlinearity would cause a singularity to form during the vortex sheet's evolution [4, 5]. An asymptotic analysis by Moore [24, 26] supports this conjecture, indicating that with initial perturbation amplitude  $\varepsilon$ , a  $\Gamma^{3/2}$  branch point forms in the vortex sheet at a finite critical time  $t = t_c(\varepsilon)$ . Meiron et al. [23] analyzed the Taylor series coefficients of  $z(\Gamma, t)$

with respect to the time variable and obtained results consistent with Moore's. The validity of Moore's approximation for  $t < t_c$  has been proven [6] and special solutions have been studied [7, 14], but proving that a singularity forms for general initial data is an open problem.

### 3. The Point Vortex Approximation

Rosenhead performed the first vortex sheet computation in 1931 [31], using the periodic Cauchy kernel in (1). The sheet was discretized by a finite number of point vortices per period  $z_j(t) \sim z(\Gamma_j, t), j = 1, \dots, N$ , leading to the ordinary differential equations,

$$\frac{dz_j}{dt} = \sum_{k \neq j} K(z_j - z_k) N^{-1}, \quad K(z) = \frac{1}{2i} \cot \pi z. \quad (2)$$

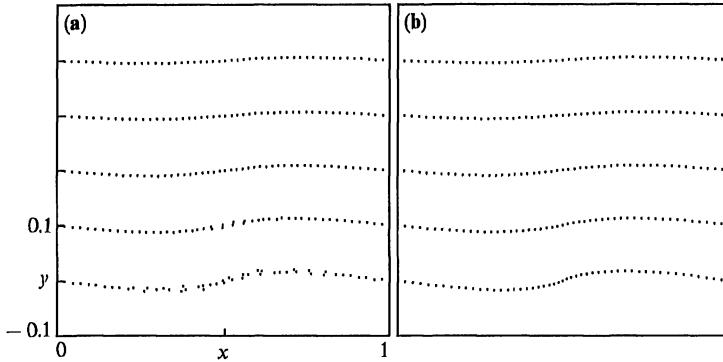
The sum omits the singular term  $k = j$ , but if the vortex sheet has a bounded 2nd  $\Gamma$ -derivative, then the discretization error is  $O(N^{-1})$  [25]. If the vortex sheet is analytic, then infinite order accuracy may be obtained by applying one step of Richardson extrapolation [34, 16].

Rosenhead used  $N \sim 10$  points and the 1st order Euler method with time step  $\Delta t \sim 0.05$  to integrate in time. He drew a smooth interpolating curve through the point vortices, suggesting that a perturbed vortex sheet rolls up into a smooth spiral. In the 1950's, Birkhoff performed computations using a larger number of point vortices and more accurate time integration [4, 5]. In contrast to Rosenhead's results, the points' computed motion was irregular, leading Birkhoff to question whether the vortex sheet rolls up into a spiral. Later workers sought to obtain convergent numerical results by using higher order accurate quadrature rules for the principal value integral, e.g. [15, 38]. Another approach was to stabilize the problem by adding surface tension [28]. In spite of much effort, the computations failed to converge as the number of points increased.

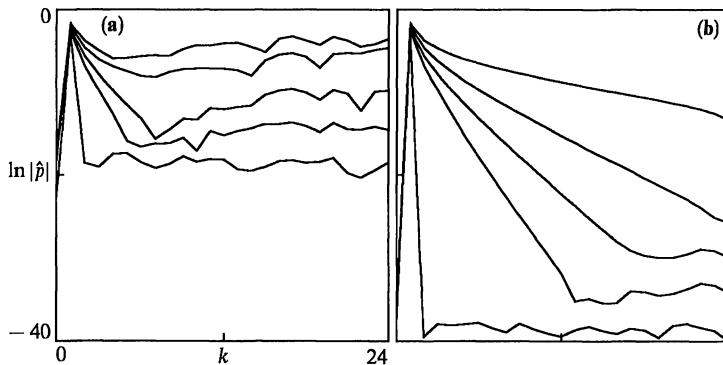
The key to obtaining convergent numerical results for  $t < t_c$  lies in Fourier analysis of the computed solution [19]. Sulem et al. [36] showed that the singularity structure of nonlinear evolution equations could be obtained from spectral computations, by analyzing the rate of decay of the discrete Fourier coefficients. For vortex sheet computations, discrete Fourier coefficients of the perturbation quantities  $p_j(t) = z_j(j) - \Gamma_j$  can be similarly analyzed.

Figure 1 shows computations with  $N = 50$  in single and double precision arithmetic. Irregular small scale motion develops in single precision, but the double precision results are smooth. The corresponding spectral amplitudes are plotted in Fig. 2. The initial spectrum has a spike at wavenumber  $k = 1$  (an explicit perturbation of amplitude  $\epsilon = 0.01$ ), as well as broad band noise in the higher modes. In Fig. 2a, the noise is amplified by the system's instability, leading to the irregular motion in Fig. 1a for  $t \geq 0.3$ . In Fig. 2b, the spectrum spreads smoothly to higher wavenumbers, due to genuine nonlinear effects [19].

A stable physical process is modeled by a well-posed initial value problem, and if the difference scheme is consistent and stable, then the solution converges as the mesh is refined [30]. Shear flows however are physically unstable and this appears as ill-posedness in the vortex sheet initial value problem. The point vortex approximation for an analytic vortex sheet defines a consistent but unstable



**Fig. 1.** Point vortex computations at times  $t = 0, 0.1, 0.2, 0.3, 0.4$ . (a) single precision.  
(b) double precision



**Fig. 2.** Discrete Fourier coefficient amplitudes corresponding to Fig. 1. (a) single precision.  
(b) double precision

difference scheme. Fritz John has observed [18], “Instability of a difference scheme under small perturbations does not exclude the possibility that in special cases the scheme converges towards the correct function, if no errors are permitted in the data or the computation.” This refers to roundoff error, due to the computer’s finite precision arithmetic, as opposed to discretization error, due to replacing a continuous operator by a discrete approximation. Using higher precision arithmetic is one way to see convergence as the mesh is refined, but for vortex sheet computations, a more practical remedy is to filter out the spurious roundoff error perturbations [19]. Computations and theory [8] now show that the point vortex approximation converges as  $N \rightarrow \infty$  for  $t < t_c$ . A consistent picture of singularity formation in a vortex sheet has been obtained: infinite curvature forms at an isolated point, but the vortex sheet remains continuously differentiable at  $t = t_c$ , showing no sign of roll-up [19, 23, 24, 26, 34].

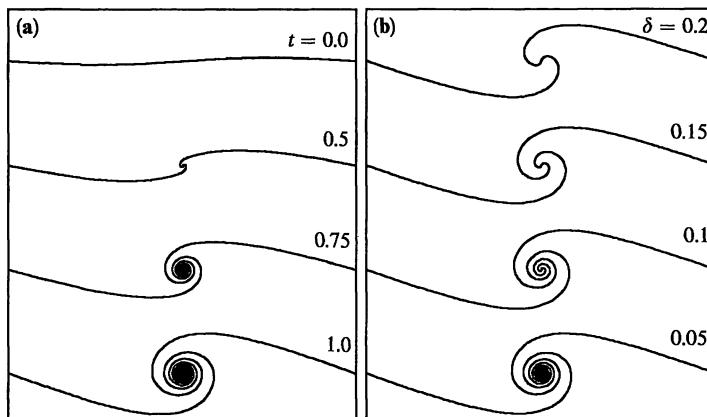
An obvious question is whether the vortex sheet continues to exist past the critical time. Note that in other problems, a physically valid weak solution can be defined past a critical time, e.g. shock formation in a nonlinear hyperbolic equation. Computations show that the point vortex approximation does not converge for  $t > t_c$  as  $N \rightarrow \infty$  [19]. A different type of small scale motion occurs in the point vortex system (2) for  $t > t_c$ , but it is not relevant to vortex sheet evolution. Based on work with self-similar vortex sheets [27], Pullin conjectured that a periodically perturbed sheet rolls up into a spiral for  $t > t_c$ , the spiral vanishes in size as  $t \rightarrow t_c^+$ , and for any  $t > t_c$  it has an infinite number of turns [29]. As described in the next section, numerical experiments using Chorin's vortex blob method support this conjecture [1, 9, 10, 11, 20, 21, 22].

#### 4. Vortex Sheet Roll-Up

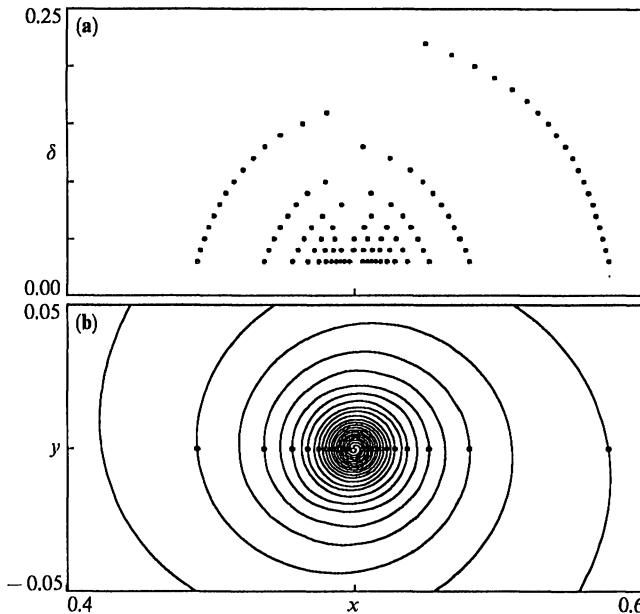
Let  $\delta > 0$  be a smoothing parameter and consider a regularized approximation to (1),

$$\frac{\overline{\partial z}}{\partial t}(\Gamma, t) = \int_{-\infty}^{\infty} K_\delta(z(\Gamma, t) - z(\tilde{\Gamma}, t)) d\tilde{\Gamma}, \quad K_\delta(z) = K(z) \frac{|z|^2}{|z|^2 + \delta^2}. \quad (3)$$

When (3) is discretized, the computational elements are called "vortex blobs". For fixed  $\delta > 0$ , short wavelength perturbations no longer have unbounded growth rates and computed solutions converge as the number of blobs  $N \rightarrow \infty$ , even for  $t > t_c$  [20]. Figure 3a shows the evolution for  $0 \leq t \leq 1$ , with the smoothing parameter value  $\delta = 0.03$ , in a case for which the vortex sheet's critical time is  $t_c \sim 0.375$ . Figure 3b shows the solution at time  $t = 1$  with decreasing amounts of smoothing  $0.05 \leq \delta \leq 0.2$ . Figure 4 shows that the smoothed solutions at  $t = 1$  converge to a spiral as  $\delta \rightarrow 0$ . The limit spiral is a candidate extension for the vortex sheet past the critical time.



**Fig. 3.** Regularized vortex sheet roll-up past the critical time  $t_c \sim 0.375$ . (a)  $\delta = 0.03$ , increasing time. (b)  $t = 1$ , decreasing  $\delta$



**Fig. 4.** Convergence as  $\delta \rightarrow 0$  for  $t = 1 > t_c \sim 0.375$ . (a) x-axis intercepts of one spiral branch plotted against  $\delta$ . (b) closeup of the solution for  $\delta = 0.03$

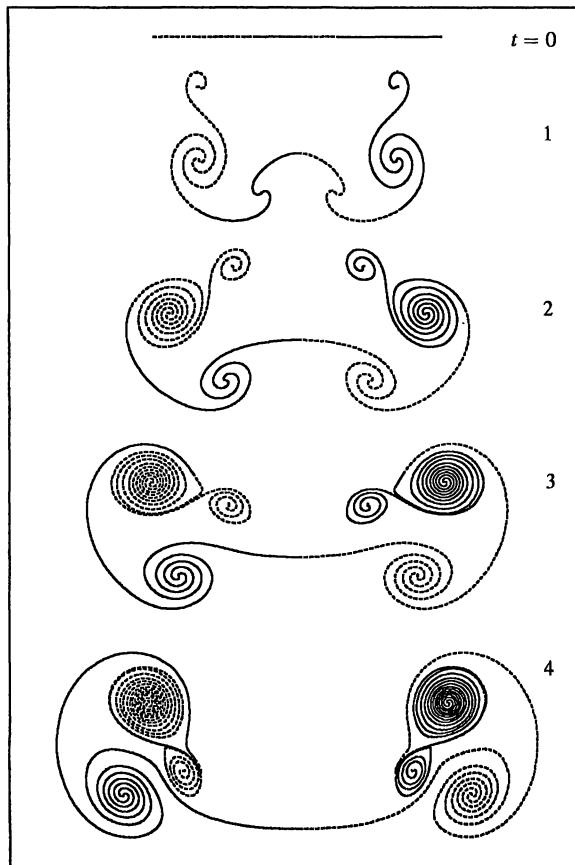
The numerical experiments suggest that the vortex blob method provides a convergent discretization of vortex sheet motion for  $t > 0$ . This has been proven for  $t < t_c$  [8], but proving convergence for the physically important roll-up regime  $t \geq t_c$  is an outstanding problem. Other interesting issues concern uniqueness of the limit for different regularizations [3, 37], existence of a weak solution to the incompressible Euler equations with general vortex sheet initial data and the possible presence of concentrations in the limit  $\delta \rightarrow 0$  [12, 13].

## 5. Applications

The vortex blob method has advantageous mathematical and numerical properties, but the smoothing parameter  $\delta$  has no precise physical meaning. One would like to know whether computations performed with a value  $\delta > 0$  approximate real fluid motion. Some applications presented below demonstrate the vortex blob method's potential for simulating shear layer dynamics.

*Aircraft Trailing Vortices.* On takeoff and landing, an aircraft sheds vortices at the wing's trailing edge. Figures 5 and 6 show a free-space vortex sheet simulation of this process, including the effects of the wing tips and deployed flaps [21]. The computation illustrates different types of vortex interactions: rotation of like-sign vortex pairs, translation of opposite-sign vortex pairs, core deformation due to collision, and vortex sheet folding.

*Separation at a Sharp Edge.* Vortices are shed from the edges of a flat plate that is moving in a viscous fluid. As the viscosity is reduced, an ideal flow emerges

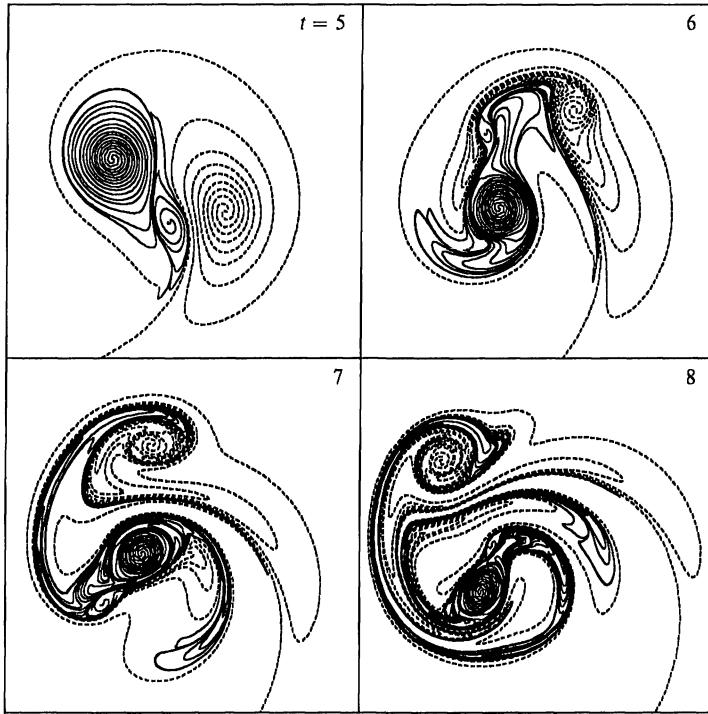


**Fig. 5.** Roll-up of an aircraft trailing vortex sheet, including tip and flap vortices. The solid and dotted lines indicate opposite senses of rotation

having embedded vortex sheets that emanate from the edges. This problem is more difficult to compute than the periodic and free space problems considered above. New issues arise, in satisfying the flow tangency condition on the plate, and shedding the correct amount of circulation at the edges. Previous numerical studies did not obtain smooth spiral roll-up, e.g. [17, 33].

A new implementation of the vortex blob method has been developed. Figure 7a is a computation of the vortex sheets that separate from an impulsively started flat plate. The velocity field plotted in Fig. 7b shows that the sheets form a recirculating region behind the plate.

To validate the algorithm, a comparison with Pullin's computation of self-similar vortex sheet roll-up [27] has been performed. The similarity assumption circumvents the difficulty of solving the initial value problem. Figure 8 compares a time dependent vortex blob computation with Pullin's self-similar result. The two plots may be superimposed to verify that the spiral shapes are in good agreement. Further details are given in [22] and a more complete validation is in preparation.



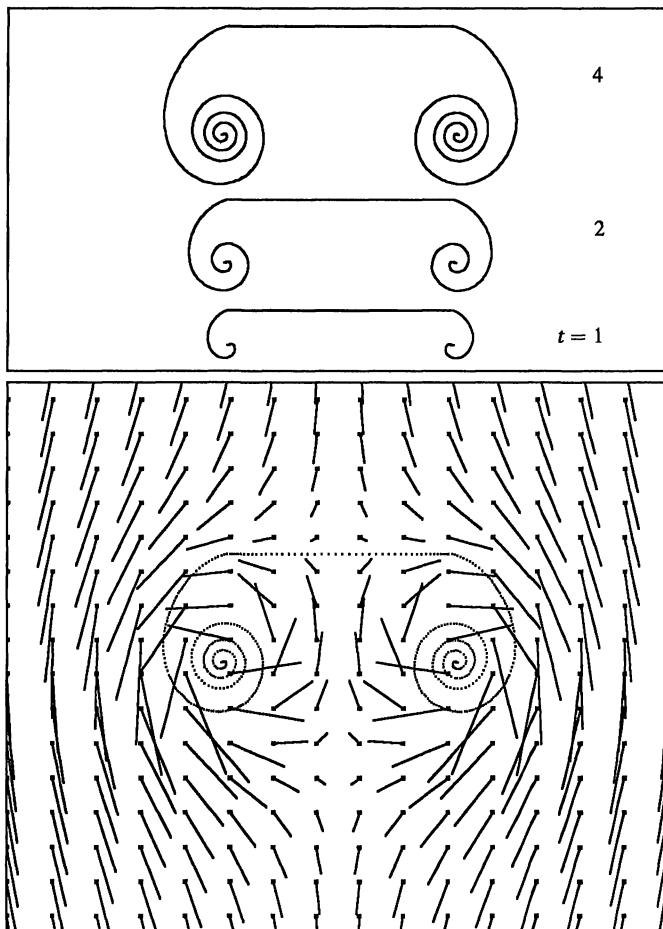
**Fig. 6.** Continuation of Fig. 5, showing details of core deformation and vortex sheet folding

*Instability of a Jet.* Figure 9 shows the evolution of a jet being expelled from a box. The jet is driven by two point sources in the lower corners of the box, which are turned on at time  $t = 0$ . A starting vortex forms and propagates away from the outlet, leaving behind a thin straight jet. Waves form along the jet, rolling up into a small vortex which propagates through the large starting vortex.

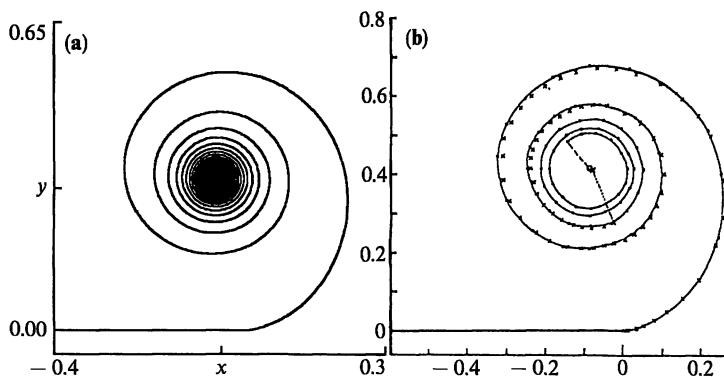
## 6. Final Remarks

Vortex sheet motion poses interesting mathematical problems concerning singular integrals, weak limits, and nonlinear dynamics. Vortex blob computations may provide a useful tool for clarifying the role of coherent vortex structures in shear flow. Future computational work will focus on improved treatment of boundary conditions, the effects of parametric forcing, and three dimensionality.

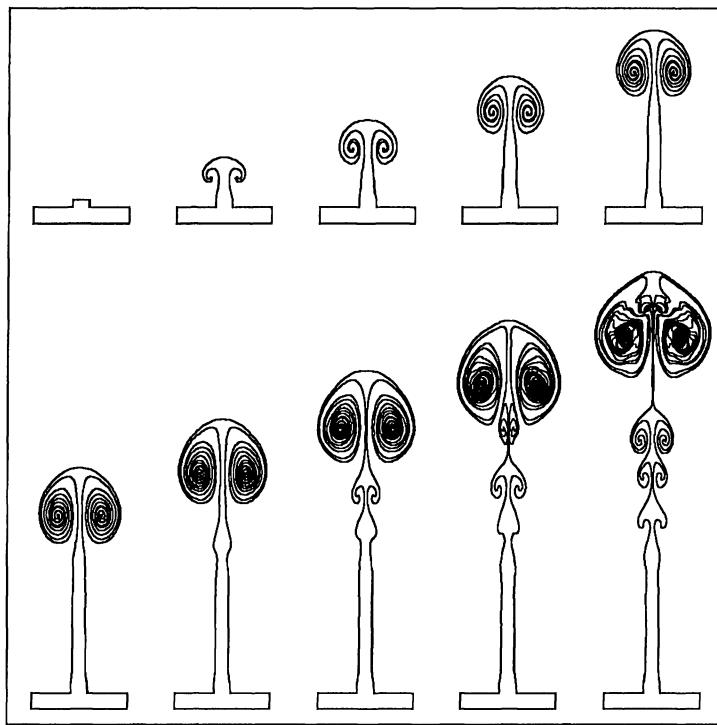
*Acknowledgements.* This work was supported in part by GRI Contract #5088-260-1692, NSF Grant DMS-#8801991, and ONR URI#N000184-86-K-0684. The computations were performed at the NSF San Diego Supercomputer Center and the University of Michigan.



**Fig. 7.** (a) Vortex sheet roll-up due to the impulsively started upward motion of a flat plate. (b) Velocity field at time  $t = 4$



**Fig. 8.** (a) Time dependent vortex blob computation,  $\delta = 0.025, t = 1$ . (b) Self-similar vortex sheet roll-up past a semi-infinite flat plate, reproduced from [27]



**Fig. 9.** Computation of a jet being expelled from a box

## References

1. Anderson, C.: A vortex method for flows with slight density variations. *J. Comp. Phys.* **61** (1985) 417
2. Anderson, C., Greengard, C. (eds.): *Vortex Dynamics and Vortex Methods*. (Lectures in Applied Mathematics, AMS). Proc. 1990 AMS-SIAM Summer Seminar in Appl. Math.
3. Baker, G.R., Shelley, M.J.: On the connection between thin vortex layers and vortex sheets. *J. Fluid Mech.* **215** (1990) 161–194
4. Birkhoff, G.: Helmholtz and Taylor instability. *Proc. Symp. Appl. Math.* **XIII** Amer. Math. Soc., Providence, R. I. (1962) 55–76
5. Birkhoff, G., Fisher, J.: Do vortex sheets roll-up? *Rend. Circ. Mat. Palermo Ser. 2*, **8** (1959) 77–90
6. Caflisch, R., Orellana, O.: Long time existence for a slightly perturbed vortex sheet. *Comm. Pure Appl. Math.* **39** (1986) 807–838
7. Caflisch, R., Orellana, O.: Singular solutions and ill-posedness for the evolution of vortex sheets. *SIAM J. Math. Anal.* **20** (1989) 293–307
8. Caflisch, R., Lowengrub, J.: Convergence of the vortex method for vortex sheets. *SIAM J. Numer. Anal.* **26** (1989) 1060–1080
9. Chorin, A.J.: Numerical study of slightly viscous flow. *J. Fluid Mech.* **57** (1973) 785–796

10. Chorin, A.J., Bernard, P.S.: Discretization of a vortex sheet with an example of roll-up. *J. Comp. Phys.* **13** (1973) 423–429
11. Chorin, A.J.: Computational fluid mechanics: Selected papers. Academic Press, Boston 1989
12. DiPerna, R.J., Majda, A.J.: Oscillations and concentrations in weak solutions of the incompressible fluid equations. *Commun. Math. Phys.* **108** (1987) 667–689
13. DiPerna, R.J., Majda, A.J.: Concentrations in regularizations for 2-D incompressible flow. *Comm. Pure Appl. Math.* **XL** (1987) 301–345
14. Duchon, J., Robert, R.: Global vortex sheet solutions of Euler equations in the plane. *J. Diff. Eq.* **73** (1988) 215–224
15. Higdon, J.J.L., Pozrikidis, C.: The self-induced motion of vortex sheets. *J. Fluid Mech.* **150** (1985) 203–231
16. Hou, T.Y., Lowengrub, J., Krasny, R.: Convergence of a point vortex method for vortex sheets. *SIAM J. Numer. Anal.* **28** (1991) 308–320
17. Graham, J.M.R.: Application of discrete vortex methods to the computation of separated flows. In: Morton, K.W., Baines, M.J. (eds.), *Numerical methods for fluid dynamics II*. Clarendon Press 1985, pp. 273–302
18. John, F.: Partial differential equations, 4th ed. Springer, Berlin Heidelberg New York 1982, p. 8
19. Krasny, R.: A study of singularity formation in a vortex sheet by the point vortex approximation. *J. Fluid Mech.* **167** (1986) 65–93
20. Krasny, R.: Desingularization of periodic vortex sheet roll-up. *J. Comp. Phys.* **65** (1986) 292–313
21. Krasny, R.: Computation of vortex sheet roll-up in the Trefftz plane. *J. Fluid Mech.* **184** (1987) 123–155
22. Krasny, R.: Vortex sheet computations: Roll-up, wakes, separation. In: Anderson, C., Greengard, C. (eds.), *Vortex dynamics and vortex methods*. (Lectures in Applied Mathematics, AMS). Proc. 1990 AMS-SIAM Summer Seminar in Appl. Math.
23. Meiron, D.I., Baker, G.R., Orszag, S. A.: Analytic structure of vortex sheet dynamics. 1. Kelvin-Helmholtz instability. *J. Fluid Mech.* **114** (1982) 283–298
24. Moore, D.W.: The spontaneous appearance of a singularity in the shape of an evolving vortex sheet. *Proc. Roy. Soc. Lond. A* **365** (1979) 105–119
25. Moore, D.W.: On the point vortex Method. *SIAM J. Sci. Stat. Comp.* **2** (1981) 65–84
26. Moore, D.W.: Numerical and analytical aspects of Helmholtz instability. In: F.I. Niordson, N. Olhoff (eds.), *Theoretical and applied mechanics*. Proc. XVI IUTAM Conf. North-Holland, Amsterdam 1984, pp. 629–633
27. Pullin, D.I.: The large-scale structure of unsteady self-similar rolled-up vortex sheets. *J. Fluid Mech.* **88** (1978) 401–430
28. Pullin, D.I.: Numerical studies of surface-tension effects in nonlinear Kelvin-Helmholtz and Rayleigh-Taylor instability. *J. Fluid Mech.* **119** (1982) 507–532
29. Pullin, D.I.: Personal communication, 1983
30. Richtmyer, R.D., Morton, K.W.: Difference methods for initial-value problems. Interscience, New York London Sydney 1967
31. Rosenhead, L.: The formation of vortices from a surface of discontinuity. *Proc. R. Soc. Lond. A* **134**, 170–192
32. Rott, N.: Diffraction of a weak shock with vortex generation. *J. Fluid Mech.* **1** (1956) 111–128
33. Sarpkaya, T.: An inviscid model of two-dimensional vortex shedding for transient and asymptotically steady separated flow over an inclined plate. *J. Fluid Mech.* **68** (1975) 109–128
34. Shelley, M.: A study of singularity formation in vortex sheet motion by a spectrally accurate vortex method. *J. Fluid Mech.* (to appear 1991)

35. Sulem, C., Sulem, P.L., Bardos, C., Frisch, U.: Finite time analyticity for the two- and three-dimensional Kelvin-Helmholtz instability. *Comm. Math. Phys.* **80** (1981) 485–516
36. Sulem, C., Sulem, P.L., Frisch, U.: Tracing complex singularities with spectral methods. *J. Comp. Phys.* **50** (1983) 138
37. Tryggvason, G.: Simulations of vortex sheet roll-up by vortex methods. *J. Comp. Phys.* **75** (1988) 253
38. van de Vooren, A.I.: A numerical investigation of the rolling up of vortex sheets. *Proc. R. Soc. Lond. A* **373** (1980) 67–91



# Developments in the Double Exponential Formulas for Numerical Integration

Masatake Mori

Department of Applied Physics, Faculty of Engineering,  
University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113, Japan

## 1. Optimality of the Trapezoidal Rule

The double exponential formula, abbreviated as the DE-formula, was first presented by Takahasi and Mori [18] in 1974 as an efficient and robust quadrature formula to compute integrals with end point singularity, e.g.

$$I = \int_{-1}^1 \frac{dx}{(x-2)(1-x)^{1/4}(1+x)^{3/4}} , \quad (1)$$

or over the half infinite interval, e.g.

$$I = \int_0^\infty e^{-x} \log x \sin x dx . \quad (2)$$

The DE-formula is based on the optimality of the trapezoidal rule over  $(-\infty, \infty)$  in the following sense. Consider the integral

$$I = \int_{-\infty}^\infty g(u) du , \quad (3)$$

where  $g(u)$  is analytic over  $(-\infty, \infty)$  and  $|g(u)|$  is integrable. We apply the trapezoidal rule, or equivalently the midpoint rule, to (3) with an equal mesh size  $h$ :

$$I_h = h \sum_{k=-\infty}^{\infty} g(kh) . \quad (4)$$

Then the error of (4) is expressed in terms of a contour integral [16]

$$\Delta I_h = \frac{1}{2\pi i} \int_{\hat{C}} \hat{\Phi}_h(w) g(w) dw , \quad (5)$$

where  $\hat{\Phi}_h(w)$  is called the characteristic function of the error and defined by

$$\hat{\Phi}_h(w) = \begin{cases} \frac{+2\pi i}{1 - \exp(-\frac{2\pi i}{h}w)} ; & \text{Im } w > 0 \\ \frac{-2\pi i}{1 - \exp(+\frac{2\pi i}{h}w)} ; & \text{Im } w < 0 \end{cases} \quad (6)$$

and  $\hat{C}$  consists of two infinite curves one of which runs to the left above the real axis and the other to the right below the real axis in such a way that there exists no singularity of  $g(w)$  between these two curves. Since  $\Delta I_A$  is a linear functional over the family of analytic functions on  $(-\infty, \infty)$ , it can be regarded as a Sato's hyperfunction [10] and  $-\hat{\Phi}_A(w)/(2\pi i)$  is nothing but its defining function.

Although there are infinite number of quadrature formulas for the integral (3), the trapezoidal rule (4) is proved to be optimal in the following sense. Let an arbitrary quadrature formula for (3) be

$$I_A = \sum_{k=-\infty}^{\infty} A_k g(a_k) . \quad (7)$$

Then its error is expressed also in terms of the contour integral

$$\Delta I_A = I - I_A = \frac{1}{2\pi i} \int_{\hat{C}} \hat{\Phi}_A(w) g(w) dw . \quad (8)$$

Since  $|\hat{\Phi}_A(w)|$  usually decays exponentially as  $|\operatorname{Im} w|$  becomes large for quadrature formulas of practical use we define the average decay rate  $r$  of  $|\hat{\Phi}_A(w)|$  for large  $|\operatorname{Im} w|$  as follows:

$$r = \lim_{d \rightarrow \infty} \left( \lim_{R \rightarrow \infty} \frac{1}{2R} \int_{-R+id}^{R+id} \left\{ -\frac{\partial}{\partial v} \log |\hat{\Phi}_A(w)| \right\} dw \right), \quad w = u + iv . \quad (9)$$

It is easy to see that the error of numerical integration is smaller if the decay rate  $r$  is larger. Then we have

**Theorem** (Takahasi-Mori, 1970). *Suppose that  $A_k$  and  $a_k$  in  $I_A$  for  $I$  satisfy*

$$\sum_{k=-\infty}^{\infty} \frac{|A_k|}{|a_k|^2} \leq M < \infty . \quad (10)$$

*Then, among quadrature formulas  $I_A$  whose average density of  $a_k$ 's per unit length is equal ( $= v_P$ ) to each other, the trapezoidal rule  $I_h$  with equal mesh size  $h = 1/v_P$  is optimal in the sense that  $r$  attains its maximum*

$$r_{\max} = 2\pi v_P = \frac{2\pi}{h} \quad (11)$$

*by the trapezoidal rule.*

The decay rate in case of the Simpson's rule is  $r = \pi v_P = \pi/h$ , so that the Simpson's rule is as twice inefficient as the trapezoidal rule.

## 2. The Double Exponential Formula

Now that the trapezoidal rule over  $(-\infty, \infty)$  is optimal, we can get a new efficient quadrature formula by means of a variable transformation. Let the given integral be

$$I = \int_a^b f(x) dx . \quad (12)$$

The variable transformation

$$x = \phi(u), \quad \phi(-\infty) = a, \quad \phi(+\infty) = b \quad (13)$$

leads to

$$I = \int_{-\infty}^{\infty} f(\phi(u)) \phi'(u) du . \quad (14)$$

Since this is an integral over  $(-\infty, \infty)$  we apply the trapezoidal rule with equal mesh size  $h$ , which results in a quadrature formula

$$I_h = h \sum_{k=-\infty}^{\infty} f(\phi(kh)) \phi'(kh) . \quad (15)$$

This is an infinite summation and in actual computation we need to truncate the sum

$$I_h^{(N)} = h \sum_{k=-N_-}^{N_+} f(\phi(kh)) \phi'(kh) , \quad (16)$$

where  $N = N_- + N_+ + 1$  is the number of function evaluations. Therefore the overall error of (16) is

$$\Delta I_h = I - I_h^{(N)} = I - I_h + I_h - I_h^{(N)} = \Delta I_h + \varepsilon_t , \quad (17)$$

where  $\Delta I_h$  is the discretization error defined by

$$\begin{aligned} \Delta I_h &= I - I_h = \int_{-\infty}^{\infty} f(\phi(u)) \phi'(u) du - h \sum_{k=-\infty}^{\infty} f(\phi(kh)) \phi'(kh) \\ &= \frac{1}{2\pi i} \int_{\hat{C}} \hat{\Phi}_h(w) f(\phi(w)) \phi'(w) dw , \end{aligned} \quad (18)$$

and  $\varepsilon_t$  is the truncation error defined by

$$\varepsilon_t = I_h - I_h^{(N)} = h \sum_{k=-\infty}^{-N_-} f(\phi(kh)) \phi'(kh) + h \sum_{k=N_+}^{\infty} f(\phi(kh)) \phi'(kh) . \quad (19)$$

In general if an analytic function  $g(w)$  decays rapidly as  $\operatorname{Re} w \rightarrow \pm\infty$ , then it grows rapidly as  $\operatorname{Im} w \rightarrow \pm\infty$ , and vice versa. Therefore  $|\Delta I_h|$  and  $|\varepsilon_t|$  cannot be made small at the same time and there should be an optimal decay rate of  $|g(w)|$  as  $\operatorname{Re} w \rightarrow \pm\infty$ .

In order to get an optimal quadrature formula Takahasi and Mori [18] investigated the efficiency of the formulas based on the three kinds of variable transformations  $x = \phi(u)$  which have the following asymptotic behaviors under the condition that  $\Delta I_h$  and  $\varepsilon_t$  are of the same order of magnitude:

$$(a) |\phi'(u)| \approx \exp(-|u|^m), m = 1, 2, \dots \quad (20)$$

$$(b) |\phi'(u)| \approx \exp(-c \exp |u|) \quad (21)$$

$$(c) |\phi'(u)| \approx \exp(-c \exp |u|^m), m = 3, 5, \dots \quad (22)$$

They found that the optimal decay of  $|g(u)|$  or  $|f(\phi(u))\phi'(u)|$  is double exponential, i.e.

$$|f(\phi(u))\phi'(u)| \sim \exp(-c \exp |u|), |u| \rightarrow \infty, \quad (23)$$

and the quadrature formula obtained based on this optimal transformation is called a double exponential formula, abbreviated as DE-formula.

Specifically, for the integral over  $(-1, 1)$

$$I = \int_{-1}^1 f(x) dx \quad (24)$$

the transformation

$$x = \tanh\left(\frac{\pi}{2} \sinh u\right) \quad (25)$$

gives a DE-formula, and for the integral

$$I = \int_0^\infty f_1(x) \exp(-x) dx \quad (26)$$

the transformation

$$x = \exp(u - \exp(-u)) \quad (27)$$

gives a DE-formula over  $(0, \infty)$ . It is also shown that the asymptotic error of the formula in terms of the mesh size  $h$  of the trapezoidal rule is expressed as

$$|\Delta I_h| \approx \exp\left(-\frac{C}{h}\right), \quad (28)$$

and that the asymptotic error in terms of the number  $N$  of the function evaluations is

$$|\Delta I_h^{(N)}| \approx \exp\left(-c \frac{N}{\log N}\right). \quad (29)$$

Before the DE-formula was developed a quadrature formula also based on variable transformation called the IMT-formula had been proposed by Iri, Moriguti and Takasawa in 1969 [3], which was characterized by the fact that the original finite interval of integration  $(0, 1)$  was transformed onto itself. The asymptotic error behavior of the formula was shown to be

$$\Delta I^{(N)} \approx \exp(-c\sqrt{N}), \quad (30)$$

from which we see that asymptotically the efficiency of the DE-formula is superior to that of the IMT-formula. Mori [6] presented a formula based on the transformation from  $(-1, 1)$  onto itself having an asymptotic error behavior

$$\Delta I^{(N)} \approx \exp\left(-c \frac{N}{(\log N)^2}\right) \quad (31)$$

which is slightly inferior to the DE-formula. In 1982 Murota and Iri [8] tried to improve the IMT-formula by means of parameter tuning and repeated application of the IMT-transformation and it turned out that, although the efficiency is improved by the repeated application of the IMT-transformation step by step, its limit does not attain the efficiency of the DE-formula as shown below:

$$\text{IMT-single} : \Delta I^{(N)} \approx \exp(-c\sqrt{N}) \quad (32)$$

$$\text{IMT-double} : \Delta I^{(N)} \approx \exp\left(-c \frac{N}{(\log N)^2}\right) \quad (33)$$

$$\text{IMT-triple} : \Delta I^{(N)} \approx \exp\left(-c \frac{N}{(\log N)(\log \log N)^2}\right) \quad (34)$$

...

$$\text{DE-formula} : \Delta I^{(N)} \approx \exp\left(-c \frac{N}{\log N}\right) . \quad (35)$$

### 3. Analysis of the DE-Formula on Function Spaces

At a research meeting held at the Research Institute for Mathematical Sciences of Kyoto University in 1985 M.Sugihara presented a detailed theoretical analysis on the optimality of the DE-formula introducing function spaces for integrands and his hand-written note on the analysis appeared in [12] in Japanese. Although he is now preparing a full paper about the details of the analysis, fragrance of his analysis will be worth while to be given here.

Basically he extended the analysis by Stenger [11] on  $H^p$  space to the analysis on spaces of functions defined not on the unit circle but directly on the real axis  $w \in (-\infty, \infty)$  in the  $w$ -plane, where  $w = u + iv$ . These spaces are characterized by the decay of their elements at large  $|\operatorname{Re} w|$ . First denote the strip domain in the  $w$ -plane

$$D(d) \stackrel{d}{=} \left\{ w \in \mathbb{C} \mid |\operatorname{Im} w| < \frac{\pi}{2}d \right\} \quad (36)$$

and define

$$\mathcal{A}(D(d)) \equiv \{\text{analytic functions on } D(d)\} . \quad (37)$$

He introduced a function space

$$\begin{aligned} & H_{\text{double}}(D(d); A, B) \quad (B < 1/d) \\ & \equiv \left\{ g \in \mathcal{A}(D(d)) \mid \sup_{w \in D(d)} \{|g(w)| \cdot |\exp(A \cosh Bw)|\} < +\infty \right\} \end{aligned} \quad (38)$$

with an appropriate norm  $\|g\|_{\text{double}}$ . In short this space is characterized by the double exponential decay of its elements at large  $|\operatorname{Re} w|$ . Consider the integral (3) and an approximation (7) to it in  $H_{\text{double}}$ . As to the norm of the error

$$E_N = \int_{-\infty}^{\infty} g(u) du - \sum_{k=1}^N A_k g(a_k) \quad (39)$$

Sugihara proved

**Theorem 1.** *In  $H_{\text{double}}(D(d); A, B)$*

$$\inf_{a_k, A_k \in \mathbb{R}} \|E_N(a_1, \dots, a_N; A_1, \dots, A_N)\| \geq C \exp\left(-c \frac{N}{\log N}\right). \quad (40)$$

On the other hand, the following theorem holds for the trapezoidal rule.

**Theorem 2.** *In  $H_{\text{double}}(D(d); A, B)$*

$$\|E_N(\text{trapezoidal rule})\| \leq C' \exp\left(-c' \frac{N}{\log N}\right). \quad (41)$$

From these two theorems we immediately see that the trapezoidal rule is optimal in  $H_{\text{double}}$ .

Next, in a similar way as  $H_{\text{double}}$ , Sugihara introduced another function space  $H_{\text{single}}$  characterized by the single exponential decay of its elements at large  $|\operatorname{Re} w|$ , and showed again the optimality of the trapezoidal rule in  $H_{\text{single}}$ . However, the inequality for the error in  $H_{\text{single}}$  corresponding to (41) is

$$\|E_N(\text{trapezoidal rule})\| \leq C' \exp\left(-\sqrt{\pi^2 d A} \sqrt{N + \frac{1}{2}}\right), \quad (42)$$

where  $A$  is some constant, so that it is clear that the double exponential transformation asymptotically leads to a more efficient quadrature formula than the single exponential one. Then it is quite natural to raise a question: is the trapezoidal rule more efficient in a space whose elements decay more rapidly than those in  $H_{\text{double}}$ ? Sugihara answered the question negatively by proving that there exists no element except zero function in such a space. Consequently we conclude that the DE-formula is asymptotically optimal.

#### 4. The DE-Formula for Slowly Decaying Oscillatory Integrals

Consider the integral

$$I = \int_0^{\infty} f(x) dx, \quad (43)$$

where

$$f(x) = f_1(x) \cos x, \quad f_1(x) = \text{algebraic function}. \quad (44)$$

In this case  $f(\phi(w))\phi'(w)$  does not belong to  $H_{\text{double}}$ , so that the DE-formula does not work well as seen from the analysis in the previous section. Toda and Ono

[19] applied the DE-formula followed by the Richardson's extrapolation method successfully to

$$I = \lim_{z \rightarrow 0} \int_0^\infty f(x) \exp(-zx) dx . \quad (45)$$

Afterwards Sugihara showed both theoretically and experimentally [13] that the Richardson's method is more efficiently applied to

$$I = \lim_{z \rightarrow 0} \int_0^\infty f(x) \exp(-zx^2) dx . \quad (46)$$

Recently Ooura and Mori [9] presented an interesting transformation which gives an efficient formula for integrals such as

$$I = \int_0^\infty f_1(x) \sin \omega x dx \quad (47)$$

$$I = \int_0^\infty f_1(x) \cos \omega x dx. \quad (48)$$

Consider the variable transformation

$$x = M\phi(t) = \frac{Mu}{1 - \exp(-K \sinh u)}, \quad M, K = \text{constant} . \quad (49)$$

Then  $\phi(u)$  satisfies  $\phi(-\infty) = 0$  and  $\phi(+\infty) = \infty$ . Moreover,

$$\lim_{u \rightarrow +\infty} \phi(u) = u \quad \text{double exponentially} \quad (50)$$

$$\lim_{u \rightarrow -\infty} \phi'(u) = 0 \quad \text{double exponentially} \quad (51)$$

hold. If we apply the variable transformation  $x = M\phi(u)$  to (47) and compute it by the trapezoidal rule, we have

$$I_h = Mh \sum_{k=-\infty}^{\infty} f_1(M\phi(kh)) \sin(\omega M\phi(kh)) \phi'(kh) . \quad (52)$$

Choose  $h$  such that  $\omega Mh = 2\pi$ , then

$$\sin(\omega M\phi(kh)) \sim \sin \omega Mkh = \sin 2\pi k = 0 , \quad (53)$$

so that we can truncate the sum (52) at some moderate value of  $k$ .

## 5. Problems Arising in Coding Automatic Integrator

When you write an automatic integrator based on the DE-formula you must be careful about the loss of significant digits which may occur when computing, say,  $(1+x)^{-3/4}$  in (1) in the close neighborhood of  $x = -1$ , and about the overflow. A device to avoid the loss of significant digits is given in [18].

Recently a useful method to avoid the overflow which may occur when computing the weights of the DE-formula was presented by Watanabe [21]. Consider again the integral

$$I = \int_{-1}^1 f(x)dx . \quad (54)$$

The weights of the DE-formula obtained by

$$x = \phi(u) = \tanh\left(\frac{\pi}{2} \sinh u\right) \quad (55)$$

are

$$A_k = \frac{\cosh kh}{\cosh^2\left(\frac{\pi}{2} \sinh kh\right)}, \quad k = 0, \pm 1, \pm 2, \dots , \quad (56)$$

and a careless coding often gives rise to the overflow when computing the denominator in (56) because it grows double exponentially as  $k$  becomes large. Watanabe found a recurrence relation

$$A_{k+1} = A_k \times r_k , \quad (57)$$

where

$$r_k = \frac{\cosh h + \sinh h \tanh kh}{(\cosh s_k + \phi(kh) \sinh s_k)^2} \quad (58)$$

and

$$s_k = \pi \sinh \frac{h}{2} \cosh((k + \frac{1}{2})h) , \quad (59)$$

and showed that the integral (54) can be computed by the following small code:

```

I=0
DO 10 k = N, 1, -1
    I = (I + f(ak) + f(-ak)) × rk-1
10  CONTINUE
I =  $\frac{\pi h}{2} (I + f(a_0))$  .

```

Although the denominator of  $r_k$  has a double exponential factor  $\sinh s_k$ , its inner exponential factor has a small coefficient  $\sinh(h/2)$ , so that the overflow in  $r_k$  will not occur until  $k$  becomes much larger than such  $k$  for which the overflow occurs in  $A_k$ .

## 6. Applications of the DE-Formula

The DE-formula is used for multiple integration. See [14] and [1], and the references therein.

The DE-formula is installed in many computer centers in Japan and is easily found in subroutine packages in the Japanese market. It is actually used in various fields of science and technology. In the references one paper is listed from each of the fields, the boundary element method [2], the surface charge method [20], filter analysis[4], and molecular chemistry [5]. Very recently also in the field of statistics it is proved to be quite efficient for numerical evaluation of risk of improved estimation [15]. For further reference see these papers and the references therein.

## References

1. V.U. Aihie, G.A. Evans: A comparison of the error function and the tanh transformation as progressive rules for double and triple singular integrals. *J. Comput. Appl. Math.* **30** (1990) 145–154
2. T. Higashimachi, N. Okamoto, Y. Ezawa, T. Aizawa, A. Ito: Interactive structural analysis system using the advanced boundary element method. *Boundary Elements — Proceedings of the Fifth International Conference*, Hiroshima, Japan, November 1983, eds. Brebbia, C. A., Futagami, T., Tanaka, M., A Computational Mechanics Center Publication. Springer, Berlin Heidelberg New York, pp. 847–856
3. M. Iri, S. Moriguti, Y. Takasawa: On a certain quadrature formula. *J. Comput. Appl. Math.* **17** (1987) 3–20 — translation of the original paper in Japanese that appeared in Kokyuroku, RIMS. Kyoto Univ. no. 91 (1970) 82–118
4. T. Kida, T. Kurogouchi: A simple numerical derivation of group delay and impulse response from the prescribed characteristic function of a reactance low-pass filter (in Japanese). *The Transactions of the Institute of Electronics and Communication Engineers* **J63** (1980) 421–428
5. T. Momose, T. Shida: Efficient formulas for molecular integrals over the Hiller-Sucher-Feinberg identity using Cartesian Gaussian functions: Towards the improvement of spin density calculation. *J. Chem. Phys.* **87** (1987) 2832–2846
6. M. Mori: An IMT-type double exponential formula for numerical integration. *Publ. RIMS, Kyoto Univ.* **14** (1978) 713–729
7. M. Mori: The double exponential formulas for numerical integration over the half infinite interval. *Numerical Mathematics Singapore 1988, International Series of Numerical Mathematics*, vol. 86. Birkhäuser, Boston 1988, pp. 367–379
8. K. Murota, M. Iri: Parameter tuning and repeated application of the IMT-type transformation in numerical quadrature. *Numer. Math.* **38** (1982) 327–363
9. T. Ooura, M. Mori: A quadrature formula for oscillatory integrals over the half infinite interval based on variable transformation (in Japanese). *Kokyuroku, RIMS Kyoto Univ.* no. 717 (1990) 68–75
10. M. Sato: Theory of hyperfunctions, *I. J. Fac. Sci. Univ. Tokyo* **8** (1959) 139–193
11. F. Stenger: Optimal convergence of minimum norm approximation on  $H^p$ . *Numer. Math.* **29** (1978) 345–362
12. M. Sugihara: On the optimality of the quadrature formula based on the DE-type variable transformation (in Japanese). *Kokyuroku, RIMS Kyoto Univ.* no. 585 (1986) 150–175

13. M. Sugihara: Methods of numerical integration of oscillatory functions by the DE-formula with the Richardson extrapolation. *J. Comput. Appl. Math.* **17** (1987) 47–68
14. M. Sugihara: Method of good matrices for multi-dimensional numerical integrations — An extension of the method of good lattice points. *J. Comput. Appl. Math.* **17** (1987) 197–213
15. N. Sugiura: Private communication
16. H. Takahasi, M. Mori: Error estimation in the numerical integration of analytic functions. *Rep. Comput. Centre Univ. Tokyo* **3** (1970) 41–108
17. H. Takahasi, M. Mori: Quadrature formulas obtained by variable transformation. *Numer. Math.* **21** (1973) 206–219
18. H. Takahasi, M. Mori: Double exponential formulas for numerical integration. *Publ. RIMS Kyoto Univ.* **9** (1974) 721–741
19. H. Toda, H. Ono: Some remarks for efficient usage of the double exponential formulas (in Japanese). *Kokyuroku, RIMS, Kyoto Univ.* no. 339 (1978) 74–109
20. Y. Uchikawa, T. Ohye, K. Gotoh: Improved surface charge method (in Japanese). *The Transactions of the Institute of Electrical Engineers of Japan A* **101** (1981) 263–270
21. T. Watanabe: On the double exponential formula for numerical integration (in Japanese). *Kakuyugokenkyu* **63** (1990) no. 5, 397–411

# Computational Complexity of Solving Real Algebraic Formulae

James Renegar

School of Operations Research and Industrial Engineering, Cornell University  
Ithaca, NY 14853-7501, USA

## 1. Introduction

We briefly survey recent computational complexity results for certain algebraic problems that are relevant to numerical analysis and mathematical programming. Topics include (i) linear programming, (ii) decision methods and quantifier elimination methods for the first order theory of the reals, (iii) solving real algebraic formulae approximately, and (iv) ill-posed problem instances.

## 2. Linear Programming

In this section we discuss complexity results for the linear programming problem

$$\begin{aligned} & \text{maximize } c^T x \\ & \text{subject to } Ax \geq b \end{aligned} \tag{2.1}$$

where  $c \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$  and  $A$  is an  $m \times n$  matrix.

In the last four years there has been a vast amount of work on “interior point” algorithms, motivated by Karmarkar’s algorithm [16]. Unlike the traditional simplex method which moves from vertex to vertex around the feasible region  $\{x; Ax \geq b\}$ , interior point methods proceed through the interior  $\{x; Ax > b\}$  of the feasible region.

Karmarkar’s algorithm is a “projective” interior point algorithm, the basic computation for each iteration being a projective transformation. In the last four years another breed of interior point algorithms has received a lot of attention, “path-following” algorithms. These are closer to traditional numerical analysis than are projective algorithms, having Newton’s method at their heart. The best upper bounds known for the complexity of linear programming are based on the analysis of particular path-following algorithms.

Following are the simple ideas behind the first path-following algorithm proven to have a polynomial-time bound. Let  $\alpha_i^T$  denote the  $i$ -th row of the constraint matrix  $A$ .

The *center* of the system of linear inequalities  $Ax \geq b$  is the point  $z$  which maximizes  $\prod_i (\alpha_i^T x - b_i)$ , viewed as a function restricted to the feasible region. The center exists and is unique if the feasible region is bounded and has non-empty interior, as we assume in what follows. Equivalently, then, the center is

the maximizer of the strictly concave function  $f(x) := \sum_i \ln(\alpha_i^T x - b_i)$ , viewed as being defined only on the interior of the feasible region.

The center has a natural physical interpretation which arises from the equations  $\nabla f(z) = 0$ . Consider each hyperplane  $\{x; \alpha_i^T x = b_i\}$  as emitting a force which acts on an arbitrary point  $x$  not in the hyperplane. The direction of the force is orthogonal to, and away from, the hyperplane, and its magnitude at  $x$  equals the reciprocal of the distance from  $x$  to the hyperplane, i.e., the force at  $x$  is simply  $\alpha_i / (\alpha_i^T x - b_i)$ . Then the center is the unique equilibrium point in the interior of the feasible region.

Now assume  $k^{(0)}$  is a known strict lower bound on the optimal objective value for (2.1). Let  $z^{(0)}$  denote the center of the extended system

$$\begin{aligned} Ax &\geq b \\ c^T x &\geq k^{(0)} \end{aligned} \tag{2.2}$$

and assume that  $x^{(0)}$  is known to be a feasible “good” approximation to  $z^{(0)}$ . We know  $x^{(0)}$  but not necessarily  $z^{(0)}$ . We want to move from  $x^{(0)}$  towards an optimal solution for (2.1). A natural way to proceed is to create a new system with center  $z^{(1)}$  closer to an optimal solution than  $z^{(0)}$ , and then move from  $x^{(0)}$  to a feasible “good” approximation  $x^{(1)}$  for  $z^{(1)}$ .

To create a new system we simply increase  $k^{(0)}$ . Of course we must be careful that  $k^{(0)}$  not be increased above the optimal objective value for (2.1). Hence, it is natural to replace  $k^{(0)}$  by  $k^{(1)} = \delta c^T x^{(0)} + (1 - \delta)k^{(0)}$  where  $0 < \delta \leq 1$ . This corresponds to bringing the hyperplane  $\{x; c^T x = k^{(0)}\}$  for (2.2) towards  $x^{(0)}$ , thus causing the equilibrium point to move to another point with better objective value.

To compute an approximation  $x^{(1)}$  for  $z^{(1)}$  we apply one iteration of Newton’s method, beginning at  $x^{(0)}$ , to the equations  $\nabla f^{(1)}(x) = 0$  where  $f^{(1)}(x) := \ln(c^T x - k^{(1)}) + \sum_i \ln(\alpha_i^T x - b_i)$ . (The special structure of  $f^{(1)}$  makes the gradient and Hessian very easy to compute).

It is easily proven that if  $\delta$  is sufficiently small and  $x^{(0)}$  is a “good” approximation to  $z^{(0)}$  then  $x^{(1)}$  obtained in this manner will be a “good” approximation to  $z^{(1)}$ . However, to establish noteworthy complexity bounds we need to prove something to the effect that  $\delta$  need not be “too” small. In [23], the author proved that any  $\delta$  satisfying  $0 < \delta \leq 1/13$  is allowable, i.e., proceeding iteratively the algorithm will then generate points  $x^{(j)}$ ,  $j = 0, 1, \dots$ , converging to an optimal solution. In practice, much larger  $\delta$  are acceptable.

The algorithm terminates with standard procedures for computing an exact optimal solution from sufficiently close approximations.

The number “1/13” arises from an analysis of Newton’s method for a carefully chosen coordinate system. The analysis can be greatly simplified by relying on work of Smale [32] as was shown by Renegar and Shub [28]. The latter paper presents a unified complexity analysis for several path-following interior point algorithms.

The best complexity bound known for linear programming is due to Vaidya [38]. He established the bound for a path-following algorithm not suggested to be practical, as it relies on fast matrix multiplication. Vaidya’s bound is  $O((m+n)^{3/2} n L^2 \log(L) \log \log(L))$  bit operations where (very roughly)  $L$  is the number of bits required to specify the particular problem instance, i.e., required to specify  $A$ ,  $b$  and  $c$ . This bound is an improvement on the earlier record

bound of Gonzaga [11] and Vaidya [37] (derived from a modification of the algorithm described above), which in turn was an improvement on the record bound of Karmarkar [16], which in turn was an improvement on the original polynomial-time bound of Khachiyan [17].

The complexity bounds for interior point algorithms are generally obtained by bounding the number of iterations required by an algorithm and multiplying the bound by the amount of work required per average iteration. The best iteration bound known for an interior point method is  $O(\sqrt{m+n} L)$  and was established by the author in [23]. The recent record complexity bounds arise from clever ways to (theoretically) reduce the amount of work required per average iteration. Much effort has been expended by researchers (including the author) to decrease the iteration bound, but to no avail.

Other seminal papers in the complexity theory of interior point methods include work by Kojima, Mizuno and Yoshise [18] (motivated by Megiddo [19]), Monteiro and Adler [21], and Ye [40]. The amount of recent literature on interior point algorithms is staggering. Relevant surveys have been written by Gonzaga [12], Goldfarb and Todd [10], and Megiddo [20].

All polynomial time algorithms for linear programming require polynomial time in the Turing machine sense, i.e., the number of bit operations is bounded by a polynomial function in the bit length  $L$  of the input. The number of arithmetic operations (over the rationals  $Q$ ) for all of the algorithms tends to infinity as  $L$  does, even when  $m$  and  $n$  are fixed. (By contrast, the number of arithmetic operations required by the simplex method can be bounded above by a function of  $m$  and  $n$  alone).

A major open question in the complexity theory of linear programming is whether or not there exist polynomial-time real number machine algorithms in the sense of Blum, Shub and Smale [2], i.e., is “uniform,” accepts arbitrary real number coefficients as inputs, and has arithmetic operation count bounded by a polynomial in  $m$  and  $n$ . Tardos [35] has made some progress on this question by devising an algorithm with arithmetic operation bound independent of the coordinates of the data vectors  $b$  and  $c$ , but dependent on the coefficients of  $A$  which she assumes to be integers.

“Experts” are divided in their opinions as to the answer of the question. If the answer was affirmative then most likely there would be practically important linear programming algorithms yet to be discovered. If the answer was negative then the complexity heirarchy for real number machines would be very different from that for Turing machines.

### 3. Decision Methods and Quantifier Elimination Methods

Now we move to a very general setting which includes many problems from numerical analysis and mathematical programming, e.g., eigenvalue problems, non-linear programming with multi-variate polynomial objective and constraint functions, sensitivity analysis problems, etc. All of these problems arise in the *classical* (by computational complexity standards) setting of the decision problem for the first order theory of the reals. We begin with a quick introduction for readers unfamiliar with this setting.

A *sentence* is an expression composed from certain ingredients. Letting  $\mathbb{R}$  denote a real-closed field, the following is an example of a sentence:

$$(\exists x_1 \in \mathbb{R}^{n_1})(\forall x_2 \in \mathbb{R}^{n_2}) \left[ (g_1(x_1, x_2) > 0) \vee (g_2(x_1, x_2) = 0) \right] \\ \wedge (g_3(x_1, x_2) \neq 0). \quad (3.1)$$

The ingredients are: vectors of variables ( $x_1$  and  $x_2$ ); the quantifiers  $\exists$  and  $\forall$ ; atomic predicates (e.g.  $g_1(x_1, x_2) > 0$ ) which are polynomial inequalities ( $>$ ,  $\geq$ ,  $=$ ,  $\neq$ ,  $<$ ,  $>$ ); and a Boolean function holding the atomic predicates ( $[B_1 \vee B_2] \wedge B_3$ ).

A sentence asserts something. The above sentence asserts that there exists  $x_1 \in \mathbb{R}^{n_1}$  such that for all  $x_2 \in \mathbb{R}^{n_2}$ , (i) either  $g_1(x_1, x_2) > 0$  or  $g_2(x_1, x_2) = 0$ , and (ii)  $g_3(x_1, x_2) \neq 0$ . Depending on the specific coefficients of the atomic predicate polynomials this assertion is true or it is false.

The set of all true sentences constitutes the first order theory of the reals. A *decision method* for the first order theory of the reals is an algorithm which, given any sentence, correctly determines if the sentence is true. Decision methods for the reals were first proven to exist by Tarski [36] who constructed one.

A sentence is a special case of a more general expression, called a *formula*. Here is an example of a formula:

$$(\exists x_1 \in \mathbb{R}^{n_1})(\forall x_2 \in \mathbb{R}^{n_2}) \left[ (g_1(z, x_1, x_2) > 0) \vee (g_2(z, x_1, x_2) = 0) \right] \\ \wedge (g_3(z, x_1, x_2) \neq 0). \quad (3.2)$$

A formula has one thing that a sentence does not, namely, a vector  $z \in \mathbb{R}^{n_0}$  of *free variables*. When specific values are substituted for the free variables, the formula becomes a sentence.

A vector  $\bar{z} \in \mathbb{R}^{n_0}$  is a *solution* for the formula if the sentence obtained by substituting  $\bar{z}$  is true.

Two formulae are *equivalent* if they have the same solutions.

A *quantifier elimination method* is an algorithm which, given any formula, computes an equivalent quantifier-free formula, i.e., for the above formula  $(\exists x_1 \in \mathbb{R}^{n_1})(\forall x_2 \in \mathbb{R}^{n_2})P(z, x_1, x_2)$  such a method would compute an equivalent formula  $Q(z)$  containing no quantified variables.

When a quantifier elimination method is applied to a sentence, it becomes a decision method. Thus, a quantifier elimination method is in some sense more general than a decision method.

Tarski [36] actually constructed a quantifier elimination method.

Many problems in numerical analysis and mathematical programming can be cast as the problem of computing a solution for a particular formula. The reader will easily verify that this can be done for the problems mentioned earlier. It can be done for many other problems as well. Of course determining if a solution for a formula exists can be done with a decision method. In the next section we discuss the complexity of approximating solutions for formulae.

Both (3.1) and (3.2) are said to be in prenex form, i.e., all quantifiers occur in front. More generally, a formula can be constructed from other formulae just as (3.1) was constructed from the atomic predicates.

We now present a brief survey of some complexity highlights for quantifier elimination methods, considering only formulae in prenex form. General bounds

follow inductively. (If a formula is constructed from other formulae, first apply quantifier elimination to the innermost formulae, then to the innermost formulae of the resulting formula, etc.)

We consider the general formula

$$(Q_1 x_1 \in \mathbb{R}^{n_1}) \dots (Q_\omega x_\omega \in \mathbb{R}^{n_\omega}) P(z, x_1, \dots, x_\omega), \quad (3.3)$$

where  $Q_1, \dots, Q_\omega$  are quantifiers, assumed without loss of generality to alternate, i.e.,  $Q_i$  is not the same as  $Q_{i+1}$ . Let  $m$  denote the number of distinct polynomials occurring among the atomic predicates and let  $d \geq 2$  be an upper bound on their degrees.

In the case of Turing machine computations where all polynomial coefficients are restricted to be integers, we let  $L$  denote the maximal bit length of the coefficients. In this context we refer to the number of “bit operations” required by a quantifier elimination method. In the general and idealized case that the coefficients are not integers we rely on the computational model of Blum, Shub and Smale [2], and refer to “arithmetic operations”, these essentially being field operations, including comparisons.

The sequential bit operation bounds that have appeared in the literature are all basically of the form

$$(md)^E [L^{O(1)} + \text{Cost}] \quad (3.4)$$

where  $E$  is some exponent and “Cost” is the worst-case cost of evaluating the Boolean function holding the atomic predicates, i.e., worst-case over 0–1 vectors.

The first reasonable upper bound for a quantifier elimination method was proven by Collins [5]. He obtained  $E = 2^{O(n)}$  where  $n := n_0 + \dots + n_\omega$ . Collins’ bound is thus “doubly exponential” in the number of variables. His method requires the formula coefficients to be integers, the number of arithmetic operations (not just bit operations) growing with the size of the integers. This is reminiscent of the polynomial time algorithms for linear programming discussed earlier. Also, Collins’ algorithm was not shown to parallelize, although enough is now known that a parallel version probably could be developed. Collins’ work has been enormously influential in the area.

The next major complexity breakthrough was made by Grigor’ev [13] who developed a decision method for which  $E \approx [O(n)]^{4\omega}$ . Grigor’ev’s bound is doubly exponential only in the number of quantifier alternations. Many interesting problems can be cast as sentences with only a few quantifier alternations. For these, Grigor’ev’s result is obviously significant. Like Collins’ quantifier elimination method, Grigor’ev’s decision method requires integer coefficients and was not proven to completely parallelize.

Slightly incomplete ideas of Ben-Or, Kozen and Reif [1] were completed by Fitchas, Galligo and Morgenstern [9] to construct a quantifier elimination method with arithmetic operation bound

$$(md)^E \text{Cost} \quad (3.5)$$

where  $E = 2^{O(n)}$ . This provides an arithmetic operation analog of Collins’ bit operation bound. When restricted to integer coefficients, the method also yields the Collins’ bound if the arithmetic operations are carried out bit by bit. Moreover, the algorithm parallelizes. Assuming each arithmetic operation requires one time unit, the resulting time bound is

$$[E \log(md)]^{O(1)} + \text{Time}(N) \quad (3.6)$$

if  $(md)^E N$  parallel processors are used, where  $\text{Time}(N)$  is the worst-case time required to evaluate the Boolean function holding the atomic predicates using  $N$  parallel processors. The analogous time bound for bit operations is also valid, namely,  $[E \log(Lmd)]^{O(1)} + \text{Time}(N)$  if  $(md)^E [L^{O(1)} + N]$  parallel processors are used.

In [25], the author introduced a new quantifier elimination method for which  $E = \prod_{k=0}^{\omega} O(n_k)$ . This was established for arithmetic operations and bit operations, i.e., (3.4) and (3.5). Regarding bit operations, the dependence of the bounds on  $L$  was shown to be very low,  $L^{O(1)}$  in (3.4) being replaced by  $L(\log L)(\log \log L)$ , i.e., the best bound known for multiplying two  $L$ -bit integers. Moreover, the method was shown to parallelize, resulting in the arithmetic operation time bound (3.6) if  $(md)^E N$  parallel processors are used, and the bit operation time bound  $\log(L)[E \log(md)]^{O(1)} + \text{Time}(N)$  if  $(md)^E [L^2 + N]$  parallel processors are used.

Independently and simultaneously, Heintz, Roy and Solerno [15] developed a quantifier elimination method for which  $E = O(n)^{O(\omega)}$ , both for arithmetic and bit operations. Their method also completely parallelizes.

The various bounds are best understood by realizing that quantifier elimination methods typically work by passing through a formula from back to front. First the vector  $x_w$  is focused on, then the vector  $x_{w-1}$ , and so on. Some methods ([5, 25]) make a second pass, from front to back. The work arising from each vector results in a factor for  $E$ . For Collins' quantifier elimination method, the factor corresponding to  $x_k$  is  $2^{O(n_k)}(2^{O(n_0)} = 2^{O(n_0)} \dots 2^{O(n_\omega)})$ . For the method introduced by the author, the factor is  $O(n_k)$ . The factor corresponding to Grigor'ev's decision method is  $\approx O(n)^4$  independently of the number of variables in  $x_k$ . In that method a vector with few variables can potentially create as much work as one with many variables. Similarly, the factor corresponding to the quantifier elimination method of Heintz, Roy and Solerno is  $O(n)^{O(1)}$  independently of  $x_k$ .

For the record, the quantifier elimination method in [25] produces a quantifier-free formula of the form

$$\bigvee_{i=1}^I \bigwedge_{j=1}^{J_i} (h_{ij}(z) \Delta_{ij} 0)$$

where  $I \leq (md)^E$ ,  $J_i \leq (md)^{E/n_0}$ ,  $E = \prod_{k=0}^{\omega} O(n_k)$ , the degree of  $h_{ij}$  is at most  $(md)^{E/n_0}$  and the  $\Delta_{ij}$  are standard relations ( $\geq, >, =, \neq, <, \leq$ ). If the coefficients of the original formula are integers of bit length at most  $L$ , the coefficients of the polynomials  $h_{ij}$  will be integers of bit length at most  $(L + n_0)(md)^{E/n_0}$ .

Results of Weispfenning [39], and Davenport and Heintz [6], show the double exponential dependence on  $\omega$  of the above bound on the degrees of the polynomials  $h_{ij}$  cannot be improved in the worst case.

Fisher and Rabin [8] proved an exponential worst-case lower bound for decision methods. However, the lower bound is exponential only in the number of quantifier alternations, and is only singly exponential in that. A tremendous gap remains between the known upper and lower bounds for decision methods.

In closing this section we mention that work of Canny ([3], [4]) has been especially influential in this area in recent years, both for the techniques he has developed and employed and for the connections he has established between the area and robotics. Work of Vorobjov [14] has also been very influential.

## 4. Solving Formulae Approximately

In this section we discuss the complexity of approximating solutions for formulae, restricting attention to the field  $\mathbb{R}$  of real numbers.

The most basic problem in this vein is that of approximating roots of univariate polynomials. Sequential bounds for this problem are numerous, for various models of computation, and have been proven over many years. Discussions have been provided by Smale ([31, 33]), and Schonhage [29].

Until very recently significant time bounds for the parallel computation of roots of univariate polynomials have been missing. However, Neff [22] has proven, in the parlance of computer science, that the problem is in NC. He has shown that all roots of a univariate polynomial can be approximated to within Euclidean distance  $\epsilon > 0$  in time  $O[\log(Ld) + \log \log(4 + \frac{1}{\epsilon})]^3$  using  $[Ld \log(2 + \frac{1}{\epsilon})]^{O(1)}$  parallel processors, where  $d \geq 3$  is the degree of the univariate polynomial and  $L$  is the maximal bit length of the coefficients, assumed to be integers. Although Neff does not present an arithmetic operation time-bound for arbitrary real number coefficients, his ideas can be extended to do so. Assuming we desire to approximate all roots lying within distance  $r$  of the origin, the resulting time bound is of the form  $[\log(d) \log \log(4 + \frac{r}{\epsilon})]^{O(1)}$  if  $[d \log(2 + \frac{r}{\epsilon})]^{O(1)}$  parallel processors are used, assuming one time unit is required per processor per arithmetic operation.

Neff's result and techniques have implications beyond the univariate setting. For example, the principal bottleneck in parallelizing Collins' quantifier elimination method has been its reliance on univariate polynomial root approximation. Neff's result removes that bottleneck.

In [26], the author reduces the problem of approximating solutions for formulae (3.3) to the problem of approximating roots for univariate polynomials. Both sequential and parallel complexity bounds for the reduction are provided. Using Neff's algorithm and ignoring the cost of evaluating the Boolean function holding the atomic predicates (which generally is a relatively negligible cost), the resulting arithmetic operation time bound is  $[E \log(md) \log \log(4 + \frac{r}{\epsilon})]^{O(1)}$  if  $[(md)^E \log(2 + \frac{r}{\epsilon})]^{O(1)}$  parallel processors are used, where as in the last section,  $E = \prod_k O(n_k)$ . For each connected component of the solution set that intersects  $\{z; \|z\| \leq r\}$ , a point within Euclidean distance  $\epsilon$  of the component is computed within this time, assuming that either (i) both  $r$  and  $\epsilon$  are input to the algorithm or (ii) only  $\epsilon$  is input and  $r$  is defined to be the infimum of distances of all solutions from the origin.

The resulting bit operation time bound is  $[E \log(Lmd) \log \log(4 + \frac{1}{\epsilon})]^{O(1)}$  if  $[L \log(2 + \frac{1}{\epsilon})]^{O(1)}(md)^E$  parallel processors are used. For each connected component of the solution set, a point within distance  $\epsilon$  of the component is computed within this time.

(The solution set consists of at most  $(md)^E$  connected components. If the coefficients of the formula are all integers of bit length at most  $L$  then each connected component of the solution set has a point within distance  $2^{\bar{L}}$  of the origin, where  $\bar{L} = L(md)^E$ .)

Sequential bounds established in [26], the best presently available, are  $(md)^E \log \log(4 + \frac{r}{\epsilon})$  arithmetic operations and  $\hat{L}^2 \log(\hat{L}) \log \log(\hat{L})(md)^E$  bit operations where  $\hat{L} = L + \log(2 + \frac{1}{\epsilon})$ . The dependence of this arithmetic operation bound on  $r$  and  $\epsilon$  cannot be improved, as was proven by the author in [24].

The occurrence of the ratio  $r/\varepsilon$  in both the sequential and parallel arithmetic operation bounds naturally leads one to suspect that analogous bounds hold for relative error, that is, for the problem of computing a point within distance  $\varepsilon/\|\bar{x}\|$  of an actual solution  $\bar{x}$ . However, it is easily proven that even for the problem of computing relative approximations to roots of quadratic polynomials, no real number model algorithm has a uniform arithmetic operation bound independent of the quadratic polynomial, depending only on the degree and relative error desired. Uniform bounds for computing points within specified relative error of a solution require that the basic algorithmic operations include more than just arithmetic operations and comparison operations, e.g., radicals.

## 5. Ill-Posed Problem Instances

Complexity theory has been developed almost exclusively for problems for which exact data is used in computations. A theory which fully incorporates the use of approximate data has yet to be developed. Central to such a development will be the notion of an ill-posed problem instance. In this section we relate an answer to the question, “Is it possible to know a problem instance is ill-posed?” A complete development can be found in [27]. Here we can only sketch a few of the ideas, rather vaguely.

Define a problem to be a formula  $P(x, y)$  with two vectors  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$  of free variables,  $\mathbb{R}$  denoting the real numbers. A *problem instance* corresponds to specifying values for  $y$ . The values are the *data* for the instance. We say “ $x$  is a *solution for instance*  $y$ ” if the pair  $(x, y)$  is a solution for the formula.

Many problems can be cast in this format. For example, linear programming with a fixed number of constraints and variables fits this format. The vector  $y$  then specifies  $A$ ,  $b$  and  $c$  in (2.1).

We assume a formula  $P(x, y)$  encoding the problem of interest is known and we assume arbitrarily accurate approximate data for the actual instance is available through an oracle. Input to the oracle is  $\delta > 0$  and output is  $\bar{y}$  strictly within error  $\delta$  of the data for the actual instance. (Very general functions are allowed in measuring solution and data errors, but for this brief synopsis assume errors are measured by norms.)

The goals:

- 1) Determine if the actual instance has a solution.
- 2) If it has a solution, compute a  $\varepsilon$ -approximate solution, i.e.,  $\bar{x} \in \mathbb{R}^n$  guaranteed to be within error  $\varepsilon$  of a solution for the actual instance.

The goals are to be achieved using approximate data and any other information available about the actual instance, e.g., it might be known the actual instance has infinitely many solutions, even if the exact data for the instance is not known. There are no restrictions on the known information, including its form, except it is required to be consistent, i.e., a contradiction cannot be deduced from it.

Depending on the known information and the exact data for the actual instance, the goals may not be achievable. Roughly, the actual instance is *ill-posed* if the goals cannot be achieved regardless of how accurate the approximate data is.

To be more definite, we introduce three definitions, which are discussed at length in [27].

**Definition.** A problem instance  $\bar{y}$  is *indistinguishable from the actual instance* if the known information regarding the actual instance does not exclude the possibility that  $\bar{y}$  is the actual instance.

**Definition.** An *acceptable algorithm* for the problem:

- 1) Accepts as input any tuple  $(\bar{y}, \delta, \varepsilon)$  where  $\delta > 0$ ,  $\varepsilon > 0$  and  $\bar{y}$  might be provided by the oracle upon input  $\delta$ .
- 2) Replies one of the following three statements:
  - (a) “All instances which
    - (i) are indistinguishable from the actual instance and
    - (ii) are strictly within error  $\delta$  of  $\bar{y}$ ,  
have solutions, and  $\bar{x}$  is strictly within error  $\varepsilon$  of a solution for all such instances,” (where  $\bar{x}$  is computed by the algorithm).
  - (b) “All instances which
    - (i) are indistinguishable from the actual instance and
    - (ii) are strictly within error  $\delta$  of  $\bar{y}$ ,  
do not have solutions.”
  - (c) “Please provide better data accuracy.”
- 3) Can be proven correct, i.e., correct in the sense that whenever it replies with statement (2a) or (2b) and the input tuple satisfies the condition that  $\bar{y}$  is strictly within error  $\delta$  of data for an instance which is indistinguishable from the actual instance then the statement replied is indeed true.

The motivation for requiring that the algorithm can be proven correct in the sense of (3) is that if the algorithm does, say, reply with statement (2a) upon some input for which  $\bar{y}$  is strictly within error  $\delta$  of the actual instance then one can prove, in terms of one’s knowledge regarding the actual instance, that the point  $\bar{x}$  is indeed within error  $\varepsilon$  of a solution for the actual instance. In other words, one can be certain in terms of what one knows about the actual instance that the algorithm will not erroneously claim a certain point to be within error  $\varepsilon$  of a solution for the actual instance when in fact it is not. *Indeed, the reader should regard the definition of an acceptable algorithm to simply formalize the requirement that one be able to trust the algorithm not to reply with an incorrect answer for the actual instance.*

**Definition.** The actual instance is *definitely well-posed* if some acceptable algorithm replies (2a) or (2b) whenever  $\bar{y}$  is strictly within error  $\delta$  of the actual instance and  $\delta > 0$  is sufficiently small, where what constitutes sufficiently small  $\delta$  may depend on  $\varepsilon$ .

We do not provide a precise definition of an ill-posed problem instance. We only assume that whatever definition is chosen, it excludes ‘definitely well-posed’ instances.

In [27], the author argues that if the information known about the actual instance can be used to deduce the actual instance is not ‘definitely well-posed’ then it can be used to deduce the actual instance is ‘definitely well-posed’, a contradiction. Consequently, if the known information is consistent, it cannot be known the actual instance is ill-posed. (Somewhat inaccurately, the result amounts to saying that if one knows arbitrarily accurate approximate data is

insufficient for the goals, then that knowledge provides sufficient information to design acceptable algorithms for which accurate approximate data is sufficient, a contradiction.) Consequently, for problems corresponding to real formulae it is impossible to know if for the instance of interest it is pointless to collect better approximate data and try new algorithms (although it is certainly possible to sometimes know it is not pointless).

The impossibility of being able to know the actual instance is ill-posed is primarily a consequence of the existence of decision methods for the first order theory of the reals. If instead, for example, solutions are required to be rational vectors, examples can be easily constructed showing it is possible to deduce from the known information that the actual instance is not ‘definitely well-posed’ without arriving at a contradiction. This is discussed at greater length in [27].

Again, much work needs to be done to develop a complexity theory which incorporates the use of approximate data. Some reflections on this have been provided by Demmel [7], Smale [34] and Shub [30]. In [27], basics for a very general theory of condition numbers are developed.

## References

1. Ben-Or, M., Kozen, D., Reif, J.: The complexity of elementary algebra and geometry. *J. Comp. System Sci.* **32** (1986) 251–264
2. Blum, L., Shub, M., Smale, S.: On a theory of computation and complexity over the real numbers: NP-completeness, recursive functions and universal machines. *Bull. Amer. Math. Soc.* **21** (1989) 1–46
3. Canny, J.: *The Complexity of Robot Motion Planning*. MIT Press, Cambridge, Mass., 1989
4. Canny, J.: Some algebraic and geometric computations in PSPACE. *Proceedings of the 20th Annual ACM Symposium on the Theory of Computing* 1988, pp. 460–467
5. Collins, G.: Quantifier elimination for real closed fields by cylindrical algebraic decomposition. (*Lectures Notes in Computer Science*, vol. 33). Springer, Berlin Heidelberg New York 1975, pp. 515–532
6. Davenport, J., Heintz, J.: Real quantifier elimination is doubly exponential. *J. Symb. Comp.* **5** (1988) 29–35
7. Demmel, J.: On condition numbers and the distance to the nearest ill-posed problem. *Num. Math.* **51** (1987) 251–289
8. Fischer, M., Rabin, M.: Super-exponential complexity of Presburger arithmetic. In: *Complexity of Computations* (SIAM-AMS Proc., 7), 1974, pp. 27–41
9. Fitchas, N., Galligo, A., Morgenstern, J.: Algorithmes rapides en sequentiel et en parallel pour l’élimination de quantificateurs en géométrie élémentaire. *Séminaire Structures Ordonnées*, U.E.R. de Math. Univ. Paris VII, 1987
10. Goldfarb, D., Todd, M.: Linear programming. In: Nemhauser, G., Rinnooy Kan, A., Todd, M. (eds.) *Handbooks in Operations Research and Management Science*, vol. I: Optimization, Chapter 2. North-Holland, Amsterdam 1989
11. Gonzaga, C.: An algorithm for solving linear programming problems in  $O(n^3L)$  operations. In: Megiddo, N. (ed.) *Progress in Mathematical Programming – Interior Point and Related Methods*, Chapter 1. Springer, Berlin Heidelberg New York 1989
12. Gonzaga, C.: Path following methods for linear programming. To appear in *SIAM Review*
13. Grigor’ev, D.Yu.: The complexity of deciding Tarski algebra. *J. Symb. Comp.* **5** (1988) 65–108

14. Grigor'ev, D.Yu., Vorobjov, N.N. Jr.: Solving systems of polynomial inequalities in subexponential time. *J. Symb. Comp.* **5** (1988) 37–64
15. Heintz, J., Roy, M.-F., Solernó, P.: Sur la complexité du principe de Tarski-Seidenberg. *Bull. Soc. Math. France* **118** (1990) 101–126
16. Karmarkar, N.: A new polynomial time algorithm for linear programming. *Combinatorica* **4** (1984) 373–395
17. Khachiyan, L.G.: Polynomial algorithms in linear programming. *USSR Computational Mathematics and Mathematical Physics* **20** (1980) 53–72
18. Kojima, M., Mizuno, S., Yoshise, A.: A polynomial-time algorithm for a class of linear complementarity problems. *Mathematical Programming* **44** (1989) 1–26
19. Megiddo, N.: Pathways to the optimal set in linear programming. In: Megiddo, N. (ed.) *Progress in Mathematical Programming – Interior Point and Related Methods*, Chapter 8. Springer, Berlin Heidelberg New York 1989
20. Megiddo, N.: On the complexity of linear programming. In: Bewley, T. (ed.) *Advances in Economic Theory*. Cambridge University Press, 1987, pp. 225–268
21. Monteiro, R.C., Adler, I.: Interior path-following primal-dual algorithms. *Mathematical Programming* **44** (1989) 27–66
22. Neff, A.C.: Specified precision polynomial root isolation is in NC. To appear in *J. Comp. System Sci.*
23. Renegar, J.: A polynomial-time algorithm, based on Newton's method, for linear programming. *Mathematical Programming* **40** (1988) 59–94
24. Renegar, J.: On the worst-case arithmetic complexity of approximating zeros of polynomials. *J. Comp.* **3** (1987) 90–113
25. Renegar, J.: On the computational complexity and geometry of the first-order theory of the reals. To appear in *J. Symb. Comp.*
26. Renegar, J.: On the computational complexity of approximating solutions for real algebraic formulae. To appear in *SIAM J. Computing*
27. Renegar, J.: Is it possible to know a problem instance is ill-posed? To appear in: *From Topology to Computation, Proceedings of the Smalefest*
28. Renegar, J., Shub, M.: Unified complexity analysis for Newton LP methods. To appear in *Mathematical Programming*
29. Schonhage, A.: Equation solving in terms of computational complexity. *Proceedings of the International Congress of Mathematicians*, Berkeley, CA, 1986. American Mathematical Society, 1987, pp. 131–153
30. Shub, M.: On the work of Steve Smale on the theory of computation. To appear in: *From Topology to Computation, Proceedings of the Smalefest*
31. Smale, S.: On the efficiency of algorithms of analysis. *Bull. Amer. Math. Soc.* **13** (1985) 87–121
32. Smale, S.: Newton's method estimates from data at one point. *The Merging of Disciplines in Pure, Applied and Computational Mathematics*. Springer, Berlin Heidelberg New York 1986, pp. 185–196
33. Smale, S.: Algorithms for solving equations. *Proceedings of the International Congress of Mathematicians*, Berkeley, CA, 1986. American Mathematical Society, 1987, pp. 87–121
34. Smale, S.: Some remarks on the foundations of numerical analysis. *SIAM Review* **32** (1990) 211–220
35. Tardos, E.: A strongly polynomial algorithm for solving combinatorial linear programs. *Operations Research* **34** (1986) 250–256
36. Tarski, A.: *A Decision Method for Elementary Algebra and Geometry*. University of California Press, 1951
37. Vaidya, P.: An algorithm for linear programming which requires  $O(((m+n)n^2 + (m+n)^{1.5}n)L)$  arithmetic operations. *Mathematical Programming* **47** (1991) 175–202

38. Vaidya, P.: Speeding-up linear programming using fast matrix multiplication. Proceedings of the 30th IEEE Symposium on the Foundations of Computer Science, pp. 332–337
39. Weispfenning, V.: The complexity of linear problems in fields. J. Symb. Comp. 5 (1988) 3–27
40. Ye, Y.: An  $O(n^3L)$  potential reduction algorithm for linear programming. To appear in Mathematical Programming

# Turbulence, Dynamical Systems and the Unreasonable Effectiveness of Empirical Eigenfunctions

Philip Holmes<sup>1</sup>, Gal Berkooz<sup>2</sup>, and John L. Lumley<sup>3</sup>

<sup>1</sup> Departments of Theoretical and Applied Mechanics and Mathematics and Center for Applied Mathematics, Cornell University, Ithaca, NY 14853, USA

<sup>2</sup> Center for Applied Mathematics, Cornell University, Ithaca, NY 14853, USA

<sup>3</sup> Sibley School of Mechanical and Aerospace Engineering and Center for Applied Mathematics, Cornell University, Ithaca, NY 14853, USA

## 1. Introduction

Paul Halmos [1981] once claimed that applied mathematics is bad mathematics. Naturally, we do not share this view (nor did Halmos, at least not entirely, in the article cited above). With the analysis of turbulence as an example, in this short paper we hope to show that, while its concerns are neither as clean nor as circumscribed, applied mathematics can be rich and fascinating and that it often interacts deeply with the best “pure” mathematics.

Two threads can be detected in the application of dynamical systems theory to turbulence. The first and better known originates with Landau [1944] and Hopf [1948], who proposed a “soluble” model equation which shared features of Navier-Stokes and could be shown to exhibit a sequence of *bifurcations* to flows of increasing complexity as a parameter (~ Reynolds number) increases. Hopf also suggested that there should be a finite dimensional attracting manifold in the infinite dimensional phase space of the evolution equations. Ruelle and Takens [1971], following an idea also proposed by Arnold, introduced the notions of generic properties and structural stability to the discussion and argued that *strange attractors*, characteristic of low dimensional systems, would more likely provide an explanation for the complex, apparently statistical motions of systems ostensibly governed by the deterministic Navier-Stokes equations than the quasi-periodic flows of Landau and Hopf. This idea led to enormous activity – experimental, analytical and numerical – especially in studies of closed fluid systems such as Bénard convection and the Taylor-Couette problem (cf. Swinney and Gollub [1981]).

*Center manifold theory* and the *unfolding* of *degenerate bifurcations* (cf. Guckenheimer and Holmes [1983], Golubitsky and Guckenheimer [1986]) have been major tools in the study of interacting instability modes. However, this local approach seems best suited to problems in hydrodynamic instability and the transition to turbulence rather than the fully developed turbulence characteristic of “open” flows. More recently, proofs of finite Hausdorff dimension for attractors of various partial differential equations (PDEs), including Navier-Stokes, and of *inertial manifolds* which globally attract all initial data exponentially fast, have helped connect these finite dimensional ideas with infinite dimensional evolution equations, cf. Constantin et al. [1989], Temam [1988].

The second thread, also traceable to a paper of Hopf [1952], cf. Foias-Prodi [1967], Foias [1973], takes a statistical viewpoint and addresses the notions of *invariant measures* and other probabilistic descriptions of turbulent fields. Of course, since the original work of Reynolds [1895], statistical descriptions of turbulence have been widely used in engineering and physics.

This paper outlines aspects of current work in which we attempt to bring these statistical and deterministic approaches together. Taking the near wall region of a fully developed turbulent boundary layer as a specific case and using a basis of empirical eigenfunctions, arrived at by statistical means, the Navier-Stokes equations are projected into a low dimensional subspace and the resulting ordinary differential equations (ODEs) studied by dynamical systems techniques. See Aubry et al. [1988, 1989, 1990], Holmes [1990], Berkooz et al. [1991] for original material and background. Here we concentrate on mathematical issues involved in the projection, truncation and modelling processes and indicate how some of the “intuitive” simplifications made in the work cited above can be justified.

In Sect. 2 we discuss the proper orthogonal decomposition, by which an “optimal” basis is generated from data ensembles. Section 3 addresses the averaging implicit in representations of boundary layer flows lacking streamwise variation and shows that the many spatial scales and modes neglected in such truncations can be rationally modelled. In Sect. 4 we draw conclusions.

We hope that this brief paper provides at least a taste of the fascinating interplay between physical foundations, experimental work, modelling and diverse types of mathematical analysis characteristic of “good” applied mathematics.

## 2. The Proper Orthogonal Decomposition and “Optimality”

Lumley [1967, 1970] first suggested the use of the *proper orthogonal* or *Karhunen-Loëve decomposition* in turbulence studies (cf. Loëve [1955]). Motivated by experimental observation of coherent structures (cf. Cantwell [1981]) in open flows such as fully developed jets, wakes, shear layers and boundary layers, he sought an unbiased method for the recognition and “extraction” of such structures as objects in space-time.

We describe the method for scalar fields; the vectorial generalization is not difficult. Suppose that  $U = \{u^i(x)|i \in I\}$  is an ensemble of realizations of turbulent field on some region  $\Omega \subseteq \mathbb{R}^n$ ; each  $u^i$  belonging to a suitable Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle$  and norm  $\|\cdot\|$ . Here  $I$  is an index set for the realizations,  $\langle \cdot \rangle$  denotes the ensemble average. For simplicity one can think of time averages in a statistically stationary flow. We seek a basis  $\Phi = \{\phi_j(x)\}_1^\infty$  for  $\mathcal{H}$  such that the ensemble averaged normalized projections  $\langle P_j(u) \rangle = \langle (u, \phi_j) \rangle / \|\phi_j\|$  onto each element in turn are maximized among all bases  $\Psi$ . The desired basis is produced by solution of the Fredholm integral equation

$$\int_{\Omega} R(x, x') \phi_j(x') dx' = \lambda_j \phi_j(x), \quad (2.1)$$

where  $R(x, x') = \langle u(x)u(x') \rangle$  is the ensemble averaged two point autocorrelation function. The basis elements  $\phi_j$  are called *empirical eigenfunctions*, since they derive from  $R(x, x')$ , itself the result of experimental observation or numerical

simulation; they are orthogonal (henceforth assumed *orthonormal*) and the eigenvalues  $\lambda_j$  correspond to the ensemble averaged kinetic energy in each “mode”,  $\phi_j$ , via the expressions

$$u = \sum_1^{\infty} a_j \phi_j, \quad \|u\|^2 = \sum_1^{\infty} |a_j|^2. \quad (2.2)$$

Also, since the  $a_j$  are uncorrelated, we have

$$\left\langle \sum_1^{\infty} |a_j|^2 \right\rangle = \sum_1^{\infty} \langle |a_j|^2 \rangle = \sum_1^{\infty} \lambda_j. \quad (2.3)$$

$\Phi$  may have other desirable properties; for example, when the tensor  $R(x, x')$  derives from measurements on an incompressible fluid, the  $\phi_j$  are divergence-free vectors: an advantage in Galerkin projection of the Navier-Stokes equations, since the projected pressure term  $(\nabla p, \phi_j)$  can be removed by integration by parts.

The basis  $\Phi$  is optimal in the following sense. Let  $\Psi = \{\psi\}_1^{\infty}$  be any other orthonormal basis, so that any field  $u \in U$  can be approximated by the  $n$ 'th order truncations  $u \approx \sum_1^n a_j \phi_j$  and  $u \approx \sum_1^n b_j \psi_j$ :

**Proposition 1.** *For each  $n \geq 1$  and any  $\Psi$ ,  $\sum_1^n \langle a_j \rangle^2 \geq \sum_1^n \langle b_j \rangle^2$ .*

This may be proved by a manipulation involving the correlation matrix  $R(x, x')$  and the fact that, if  $K$  is a self-adjoint operator and  $Q$  an orthogonal projector onto  $\text{span}\{\phi_1, \dots, \phi_n\}$ , then

$$\text{Tr}(K \cdot Q) = \sum_1^n (K \phi_j, \phi_j) \leq \sum_1^n \kappa_j, \quad (2.4)$$

where  $\kappa_1, \dots, \kappa_n$  are the  $n$  largest eigenvalues of  $K$  (cf. Temam [1988, p. 260]). Proposition 1 guarantees that use of  $\Phi$  minimizes the error, in a mean square sense, among all possible truncations of any fixed order. In fact  $\Phi$  spans the subspace which contains almost all realizations in a measure theoretic sense of the flow from which  $R(x, x')$  was computed.

Several groups have recently been using empirical eigenfunctions for the representation of turbulent fields: see Moin et al. [1984, 1989] and Sirovich et al. [1987, 1988, 1989, 1990] for examples based on computer simulation. However, in spite of all the numerical activity and studies of convergence of averages (cf. Foias et al. [1990]), there has been little study of the way in which the mean square optimality of  $\Phi$  relates to the dynamics of ordinary differential systems produced via projection of the governing PDEs onto (low dimensional) subspaces spanned by finite sets  $\Phi^N = \{\phi_j\}^N$ . The original work of our group (Aubry et al. [1988]) indicated that the statistical optimality of  $\Phi^N$  led to systems which exhibited *instantaneous dynamical behavior* representative of the full system, in spite of the low order truncations employed (only 5 complex modes!). Sirovich and his colleagues have made similar observations for other dissipative PDEs. The ideas introduced in the rest of this paper are directed toward a rational justification of these observations.

### 3. Dynamics in Subspaces Lacking Streamwise Variation

In any finite dimensional representation of an infinite dimensional evolution equation, such as Navier-Stokes, the phase space  $\mathcal{H}$  is divided into a finite subset of resolved modes,  $R$ , and its (orthogonal) complement,  $S$ , elements of which are modelled in terms of elements of  $R$ , or simply neglected. For example, if an inertial manifold,  $M$ , exists, it is expressible as a function  $h : R \rightarrow S$  and one relies on the attractivity of  $M$  to guarantee that all states asymptotically approach points of  $M$  with coordinates  $(r, h(r))$ . The reduced dynamical system or *inertial form* (Sell [1989]) is then a closed set of ODEs for  $r$ : in its simplest form it is merely the projection of the full flow onto  $R$  with  $h \equiv 0$ .

In measure theoretic terms, there is a conditional measure  $\mu_r(s)$  on each unresolved fiber  $r + S$  over the base  $R \ni r$  and an associated measure  $\mu_R(r)$  on  $R$  itself, forming a measurable partition, so that for a set  $B \subset \mathcal{H}$ :

$$\mu(B) = \int_R \mu_r(s \in B) d\mu_R(r). \quad (3.1)$$

In this context, the physical notion of small scale or local isotropy (Tennekes and Lumley [1974]) is the assumption that  $\mu_r$  is independent of the base point  $r$ . More generally, modelling of activity in small spatial scales in terms of the large scales can be seen as an attempt to estimate  $\mu_r$ .

Throughout the analysis to follow two notions of ensemble average are implicit. The first is the usual one of averaging over many separate (experimental) realizations. The second is based on the observation that, when length scales in the homogeneous directions (those lacking a distinguished origin) are long compared to those of typical turbulent phenomenon, integration over those directions will yield a characteristic measure of all the dynamical phenomena. An implicit equivalence or ergodicity assumption is thus invoked.

We now focus on a turbulent boundary layer over a flat plate, the domain  $\Omega = [0, L_1] \times [0, L_2] \times [0, L_3]$  being of streamwise extent  $L_1$ , spanwise  $L_3$  and normal  $L_2$ , with periodic boundary conditions in  $x_1$  and  $x_3$ . In Aubry et al. [1988, 1989, 1990]  $\Omega$  is the “wall region”.

Flow visualizations of the boundary layer by Kline et al. [1967] (one of which is reproduced in Aubry et al. [1988]) demonstrate the presence of coherent structures (streaks) with long streamwise spatial scales and relatively small spanwise spacing. This, together with explicit evaluation of empirical eigenfunctions, prompts the split of  $\mathcal{H}$  into  $R$  and  $S$  developed below.

In this section we identify  $R$  with (a finite dimensional subset of) the subspace of flows  $u(x, t) = u(x_2, x_3, t)$  having no streamwise dependence. Thus, if  $P$  is the orthogonal projector  $P : \mathcal{H} \rightarrow R$ , its application is identical to averaging in the streamwise direction, as can be seen by appeal to the representation

$$\mathbf{u}(\mathbf{x}, t) = \sum_{k,l,n} a_{k,l,n}(t) \Phi_{k,l}^n(x), \quad (3.2)$$

where

$$\Phi_{k,l}^n(x) = e^{2\pi i \left( (kx_1/L_1) + (lx_3/L_3) \right)} \phi_{k,l}^n(x_2) / \sqrt{L_1 L_3},$$

and

$$P(u) = \sum_{l,n} a_{0,l}^n(t) \Phi_{0,l}^n(x) \equiv \frac{1}{L_1} \sum_{l,n} a_{k,l}^n(t) \int_0^{L_1} \Phi_{k,l}^n(x) dx_1. \quad (3.3)$$

Here  $u(x, t)$  is the fluctuating field riding on the mean flow  $U = (U_1(x_2), 0, 0)$ .

Let

$$\frac{du}{dt} = N(u) \quad (3.4)$$

denote the Navier-Stokes equations. We wish to determine the evolution of the resolved state  $r \in R$ . Ideally, we should solve (3.4) for an ensemble of initial conditions  $(r, s_i)$  to find  $\mu_r$  and determine the vector field in  $R$  at each such  $r$  by integration with respect to this measure:

$$\left( \frac{dr}{dt} \right)_a = \int_{S_r} P \left( \frac{du}{dt} \Big|_{u=r+s} \right) d\mu_r(s). \quad (3.5)$$

A simpler alternative is to project (3.4) onto  $R$  and solve the resulting reduced equation

$$\left( \frac{dr}{dt} \right)_b = P(N(r + s)), \quad (3.6)$$

with  $s = h(r)$  modelled in some way. This latter is computationally accessible and fortunately we have

**Proposition 2.** *For statistically stationary flows, as  $L_1 \rightarrow \infty$  so  $\left( \frac{dr}{dt} \right)_b \rightarrow \left( \frac{dr}{dt} \right)_a$ .*

*Proof.* The right hand side of (3.5) is the conditional ensemble average over  $S_r$  such that  $P(u) = r$ . However, if  $L_1$  is large enough  $\Omega$  may be divided into  $M = L_1/d \gg 1$  regions of length  $d$  in each of which the flow is statistically independent. (Thus  $d$  is assumed to be much greater than the streamwise length scale.) The ensemble average may then be written

$$\begin{aligned} \int_{S_r} P \left( \frac{du}{dt} \Big|_{u=r+s} \right) d\mu_r(s) &= \frac{1}{M} \sum_1^M \frac{dr}{dt} \Big|_{r+s} \rightarrow \frac{1}{L_1} \int_0^{L_1} \frac{dr}{dt} \Big|_{r+s} dx_1 \\ &= P(N(r + s)). \end{aligned} \quad \square$$

Physically, we assume that  $\Omega$  is long enough to contain “something of everything” at any instant, and so effectively to yield the measure  $\mu_r$ . (It is also related to the small scale isotropy assumption referred to earlier). Thus the projected evolution equations (3.6) differ from the “ideal” case (3.5) only in that the unresolved modes,  $s$ , must be modelled. We now turn to this aspect.

The evolution equation for the fluctuations  $u$ , including Reynolds stress terms due to the fact that the mean flow  $U(x) = (U_1(x_2), 0, 0)$  only solves the Navier-Stokes in a suitably averaged sense, may be written as in Aubry et al. [1988], using the Einstein summation convention:

$$\frac{\partial u}{\partial t} + u_{i,1} U_1 + U_{1,2} u_2 \delta_{i1} + u_j u_{i,j} - \langle u_{i,j} u_j \rangle = -\frac{1}{\varrho} \pi_{i,j} + \nu u_{i,j,j}. \quad (3.7)$$

Here  $\langle \cdot \rangle$  denotes the streamwise-spanwise spatial average:

$$\langle \cdot \rangle = \frac{1}{L_1 L_3} \int_0^{L_1} \int_0^{L_3} (\cdot) dx_1 dx_3,$$

$\pi$  the fluctuating pressure and  $\nu$  is the viscosity. Letting  $u_i = r_i + s_i$  and applying  $P$ , this yields for (3.6)

$$\begin{aligned} \frac{\partial r_i}{\partial t} + r_{i,1} U_1 + U_{1,2} r_2 \delta_{i1} + P(s_j s_{i,j}) - \langle s_j s_{i,j} \rangle \\ + r_j r_{i,j} - \langle r_j, r_{i,j} \rangle = -\frac{1}{\varrho} p_{,i} + \nu r_{i,j,j}, \end{aligned} \quad (3.8)$$

where  $\pi = p + q$  and  $p$  corresponds to the resolved pressure modes. Note that, since  $r_i$  is independent of  $x_1$  and  $s_i$  has zero mean in  $x_1$ , several “mixed” terms vanish: specifically, the Leonard stresses are zero.

$$P \left( (r_j s_{i,j}) + s_j r_{i,j} - \langle (r_j s_{i,j}) + (s_j r_{i,j}) \rangle \right) = 0, \quad (3.9)$$

as Aubry et al. [1988] assumed.

Now (3.8) is “closed” apart from the term  $P(s_j, s_{i,j}) - \langle s_j s_{i,j} \rangle$  which represents interaction between modes in  $R$  and  $S$ . Normally one might expect this to result in transfer of energy from  $R$  to  $S$  and from  $S$  to  $R$ , as well as modulate its transport among the  $R$  modes. However, when  $S$  represents a subspace of no streamwise dependence, and is spanned by the leading empirical eigenfunctions, we have

**Proposition 3.** (1)  $P(s_j s_{i,j}) - \langle s_j s_{i,j} \rangle$  on the average can only transfer energy from  $R$  to  $S$  or move it around in  $R$ . No energy on the average enters  $R$  from  $S$ .

(2) The ratio  $\langle r_1 r_2 \rangle / \|r\|^2$  of Reynolds stress to turbulent kinetic energy is restricted to an interval appropriate for the ensemble  $U$ . For the truncations of Aubry et al. [1988] this interval is bounded away from zero.

(3) Provided only low wavenumbers are retained in  $R$ , the ratio of energy loss from  $R$  to  $S$  to turbulent kinetic energy in  $R$  is compatible with simple Heisenberg or eddy viscosity modelling.

*Proof.* Multiplying (3.8) by  $r_i$  and averaging by  $\langle \cdot \rangle$  yields the evolution equation for the resolved turbulent kinetic energy  $\frac{\langle r_i r_i \rangle}{2}$ :

$$\begin{aligned} \frac{D}{Dt} \left\langle \frac{r_i r_i}{2} \right\rangle = -U_{1,2} \langle r_1 r_2 \rangle - \left\langle \left( \frac{r_i r_i}{2} + \frac{p}{\varrho} \right) r_j - \nu \left( \frac{r_i r_i}{2} \right)_{,j} - \nu (r_i r_j)_{,i} \right\rangle_j \\ - \langle P(s_j s_{i,j}) r_i \rangle - 2\nu \langle \varrho_{ij} \varrho_{ij} \rangle, \end{aligned} \quad (3.10)$$

where  $\varrho_{ij} = (r_{i,j} + r_{j,i})/2$  and  $D/Dt = \frac{\partial}{\partial t} + r_i \frac{\partial}{\partial x_i}$  denotes the convective derivative. In (3.10) the term  $\langle P(s_j, s_{i,j}) r_i \rangle$  alone represents interactions between  $R$  and  $S$ . Using incompressibility and the fact that  $r_i$  is independent of  $x_1$  it may be rewritten explicitly as

$$\begin{aligned}
& \left\langle \frac{1}{L_1} \int_0^{L_1} s_j s_{i,j} dx_1 r_i \right\rangle = \left\langle \frac{1}{L_1} \int_0^{L_1} (s_i s_j)_{,j} dx_1 r_i \right\rangle \\
&= \left\langle \frac{1}{L_1} s_i s_1 \Big|_0^{L_1} r_i + \frac{1}{L_1} \int_0^{L_1} (s_i s_2)_{,2} r_i dx_1 + \frac{1}{L_1} \int_0^{L_1} (s_i s_3)_{,3} r_i dx_1 \right\rangle \\
&= \left\langle \left( \frac{1}{L_1} \int_0^{L_1} s_i s_2 r_i dx_1 \right)_{,2} + \left( \frac{1}{L_1} \int_0^{L_1} s_i s_3 r_i dx_1 \right)_{,3} \right. \\
&\quad \left. - \frac{1}{L_1} \int_0^{L_1} s_i s_2 dx_1 r_{i,2} - \frac{1}{L_1} \int_0^{L_1} s_i s_3 dx_1 r_{i,3} \right\rangle \\
&= \langle P(s_i s_j r_i)_{,j} \rangle - \langle P(s_i s_j) r_{i,j} \rangle = \langle s_i s_j r_i \rangle_{,j} - \langle P(s_i s_j) r_{i,j} \rangle. \tag{3.11}
\end{aligned}$$

The first term in (3.11) represents transport of energy among modes in  $R$  while the second represents straining of modes in  $S$  by those in  $R$  – i.e. losses from  $R$  to  $S$ . No term exists for transfer of energy from  $S$  to  $R$  on the average. This establishes (1).

To establish (2) we consider the ratio  $\tau = \langle u_1 u_2 \rangle / \langle u_i u_i \rangle$ . Since this argument has already appeared (Berkooz et al. [1991]), we merely sketch it. Expanding  $\tau$  by (3.2)-(3.2) with  $k = 0$  and maximizing and minimizing over the available modes in any specified truncation  $|l| \leq L$ ,  $1 \leq n \leq N$  leads to a study of expressions of the form

$$\frac{(\phi_{0,l}^n)_1 (\phi_{0,l}^{n*})_2}{(\phi_{0,l}^n)_i (\phi_{0,l}^n)_i} \tag{3.12}$$

and thus the relative signs and magnitudes of the streamwise  $(\cdot)_1$  and spanwise  $(\cdot)_2$  components of the empirical basis vectors  $\phi_{0,l}^n$  determine the upper and lower bounds for the range of ratios which can be represented by velocity fields belonging to  $R$ . Reference to Fig. 4 of Aubry et al. [1988] shows that the lower bound is strictly positive, establishing (2).

Having shown in (1) that energy can be lost from  $R$  to  $S$  we now wish to model this effect by expressing the loss as a function purely of the resolved modes  $r_i$ . To do this we estimate the ratio

$$\begin{aligned}
\frac{\langle P(s_i s_j) r_{i,j} \rangle}{\langle \varrho_{ij} \varrho_{ij} \rangle} &\sim \frac{\langle r_{i,j} r_{i,j} \rangle^{1/2} \langle P(s_i s_j) P(s_i s_j) \rangle^{1/2}}{\langle (r_{i,j} + r_{j,i})(r_{i,j} + r_{j,i}) \rangle} \\
&\sim \frac{\langle r_{i,j} r_{i,j} \rangle^{1/2} \langle (s_i, s_i) \rangle}{\langle r_{i,j} r_{i,j} \rangle} \sim \frac{\langle (s_i, s_i) \rangle}{\langle r_{i,j} r_{i,j} \rangle^{1/2}}. \tag{3.13}
\end{aligned}$$

Here  $\sim$  means “equal within an order 1 number” (the correlation coefficient).

In the conventional Heisenberg type model as used by Aubry et al. [1988] (cf. Tennekes and Lumley [1972]) an effective (eddy) viscosity

$$v_T = \frac{\int_0^{L_2} \langle s_i s_i \rangle dx_2}{\left( L_2 \int_0^{L_2} \langle r_{i,j} r_{i,j} \rangle dx_2 \right)^{1/2}} \tag{3.14}$$

is introduced by averaging the numerator and square of the denominator of (3.13) over the wall region normal to the wall. From (3.13) a more natural choice might seem to be

$$\tilde{v}_T = \frac{1}{L_2} \int_0^{L_2} \left( \frac{\langle s_i s_j \rangle}{\langle r_{i,j} r_{i,j} \rangle^{1/2}} \right) dx_2, \quad (3.15)$$

but using the facts that the velocities  $s_i$ ,  $r_i$  and their derivatives are represented by the empirical eigenfunctions; which may in turn be well fitted by (low) order polynomials in the wall region we find that (3.14) and (3.15) are of the same order. Specifically, let

$$\langle s_i s_j \rangle \sim x_2^{2p}, \quad \langle r_{i,j} r_{i,j} \rangle \sim x_2^{2q}, \quad (3.16)$$

so that a typical scale in the unresolved mode is represented by a polynomial of order  $p$  and in the resolved modes by order  $q$  ( $p > q$ ). Then (3.14) and (3.15) yield, respectively

$$v_T \sim \frac{\sqrt{2q+1}}{2p+1} L_2^{2p-q}, \quad \tilde{v}_T \sim \frac{L_2^{2p-q}}{2p-q+1}, \quad (3.17)$$

so that

$$\frac{\tilde{v}_T}{v_T} \sim \frac{2p+1}{\sqrt{2q+1}(2p-q+1)}. \quad (3.18)$$

If, as in Aubry et al. [1988], the energy is assumed to be lost to the next higher wavenumbers (the lowest in  $S$ ), then  $p = q + 1$  and

$$\frac{\tilde{v}_T}{v_T} \sim \frac{2q+3}{(q+3)\sqrt{2q+1}} \in (c/\sqrt{q}, 1), \quad (3.19)$$

while if we assume that energy is transferred to high wavenumbers in  $S$ ,  $p \gg q$  and

$$\frac{\tilde{v}_T}{v_T} \sim \frac{1}{\sqrt{q}}. \quad (3.20)$$

In either case, if a small range of wavenumbers  $0 \leq q \leq Q$  are retained in  $R$ , as in Aubry et al. [1988] the ratio  $\tilde{v}_T/v_T$  does not vary radically and (3) is established.  $\square$

As Berkooz et al. [1991] point out, the restriction, by the projection, of the ratio  $\tau$  of Reynolds stress to kinetic energy to a range appropriate to the experimental observations is crucial to the success of low dimensional representations lacking streamwise variation in producing relevant dynamics. If  $\tau$  could drop to zero, decoupling spanwise from streamwise motions, then one expects turbulent kinetic energy to decay (cf. Moffatt [1990]). When  $\tau$  is strictly positive the turbulent fluctuations can (indeed must) extract energy from the mean velocity gradient  $U_{1,2}$  via the third term in (3.7) (the second is absent if there is no streamwise dependence).

The first conclusion of Proposition 3 is striking in that it shows that motions contained high spanwise and wall-normal wavenumbers but lacking streamwise variations (in  $R$ ) cannot on the average extract energy even from *low* streamwise wavenumber modes (in  $S$ ). This is distinct, of course, from their ability to extract energy from the mean velocity gradient referred to above.

## 4. Conclusions

We conclude that projections on the subspace  $S$ , lacking streamwise variation (and finite subspaces of it) have some special properties. These become apparent due to the convenient interpretation of the projection as an average in the streamwise direction (Proposition 2). The notion of equivalence of ensemble averages is used in an implicit way and must be considered an assumption. We then can make the following observations. The projected system is in effect a spatially averaged system, constraining the reduced dynamics to what will be physically observable. Assumptions previously made neglecting the Leonard stresses prove to be exact in an average sense. The expression used for the effective (Heisenberg) viscosity, which was based on physical intuition, is proven to be correct (within an order 1 number) in an average sense. We also observe that on the average energy does not pass from  $S$  to its complement,  $R$ , justifying our intuitive feeling that  $R$  is a fundamentally important subspace. (It was previously referred to as a “backbone” for the analysis Holmes [1990].) We also recalled a previous observation that the turbulent energy production is held in an experimentally appropriate range, again confirming our physical intuition.

In work of this type we use a wide range of mathematical tools. Statistical methods are used to extract key features from experimental data or simulations. Recent ideas from PDE, analysis and dynamical systems theory motivate our derivation of reduced (projected) dynamical systems. This is not a bag of unrelated tricks which happen to work; there are deep relations among the different pieces of mathematics and the physical problems which have prompted their development and use. Although we do not describe them here, ideas from the global theory of dynamical systems permit us to give fairly complete analyses of large ( $O(10-50)$ ) sets of ODEs, to understand the effects of symmetries on them (cf. Armbruster et al. [1988, 1989]) and the influence of noise and other perturbations on the heteroclinic attractors they possess (Stone and Holmes [1989, 1990]). Physical intuition and reasoning come to our aid when rigorous mathematical arguments are inadequate. Perhaps more significantly for mathematics, attempts to analyze problems such as turbulence continue to provide a wealth of challenging mathematical problems and even to suggest whole new fields of study.

*Acknowledgements.* The work described in this paper was supported by the Air Force Office of Scientific Research and the National Science Foundation under AFOSR 89 0226A (Wall Layers).

## References

- Armbruster, D., Guckenheimer, J., Holmes, P. (1988): Heteroclinic cycles and modulated travelling waves in systems with  $O(2)$  symmetry. *Physica* **29D**, 257–282
- Armbruster, D., Guckenheimer, J., Holmes, P. (1989): Kuramoto-Sivashinsky dynamics on the center-unstable manifold. *S.I.A.M. J. Appl. Math.* **49**, 676–691
- Aubry, N., Holmes, P., Lumley, J., Stone, E. (1988): The dynamics of coherent structures in the wall region of a turbulent boundary layer. *J. Fluid Mech.* **192**, 115–173
- Aubry, N., Lumley, J., Holmes, P.J. (1990): The effect of modelled drag reduction on the wall region. *Theor. Comp. Fluid Dyn.* **1**, 229–248
- Aubry, N., Sanghi, S. (1989): Streamwise and spanwise dynamics of the turbulent wall layer. In: *Forum on Chaotic Flow*, Ghia (ed.). Proc. ASME, New York

- Berkooz, G., Holmes, P.J., Lumley, J.L. (1991): Intermittent dynamics in simple models of the turbulent wall layer. *J. Fluid Mech.* (in press)
- Blackwelder, R.F. (1989): Some ideas on the control of near wall eddies. AIAA paper #89-1009, AIAA 2nd Shear Flow Conference
- Cantwell, B.J. (1981): Organized motion in turbulent flow. *Ann. Rev. Fluid Mech.* **13**, 457–517
- Constantin, P., Foias, C., Temam, R., Nicolaenko, B. (1989): Integral manifolds and inertial manifolds for dissipative partial differential equations. Springer, New York
- Foias, C. (1973): Statistical study of Navier-Stokes equations I and II. *Rend. Sem. Mat. Univ. Padova* **48**, 219–348 and **49** 9–123
- Foias, C., Manley, O., Sirovich, L. (1990): Empirical and Stokes eigenfunctions and the far-dissipative turbulent spectrum. *Phys. Fluids A* **2**, 464–467
- Foias, C., Prodi, G. (1967): Sur le comportement global des solutions non-stationnaires des équations de Navier-Stokes en dimension 2. *Rend Sem. Mat. Univ. Padova* **39**, 1–34
- Golubitsky, M., Guckenheimer, J. (1986) (eds.): Multiparameter bifurcation theory. *Contemp. Math.* **56**. AMS, Providence, R.I.
- Guckenheimer, J., Holmes, P. (1983): Nonlinear oscillations, dynamical systems and bifurcations of vector fields. (Applied Mathematical Sciences, vol. 42.) Springer, New York, 1986 (Second printing)
- Halmos, P. (1981): Applied mathematics is bad mathematics. In: *Mathematics tomorrow* (ed. L.A. Steen). Springer, New York, pp. 9–20
- Holmes, P.J. (1990): Can dynamical systems approach turbulence? In: *Whither turbulence: Turbulence at the crossroads* (ed. J.L. Lumley). (Lecture Notes in Physics, vol. 357). Springer, New York, pp. 195–249, 306–309
- Hopf, E. (1948): A mathematical example displaying the features of turbulence. *Commun. Pure Appl. Math.* **1**, 303–322
- Hopf, E. (1952): Statistical hydromechanics and functional calculus. *J. Rat. Mech. Anal.* **1**, 87–123
- Landau, L. (1944): On the problem of turbulence. *Dokl. Akad. Nauk. SSSR* **44**, 339–342
- Loève, M. (1955): Probability theory. Van Nostrand
- Lorenz, E.N. (1963): Deterministic nonperiodic flow. *J. Atmos. Sci.* **20**, 130–141
- Lumley, J.L. (1967): The structure of inhomogeneous turbulent flows. In: *Atmospheric turbulence and radio wave propagation* (eds. A.M. Yaglom and V.I. Tatarski). Nauka, Moscow, pp. 166–178
- Lumley, J.L. (1970): Stochastic tools in turbulence. Academic Press, New York
- Moffatt, H.K. (1990): Fixed points of turbulent dynamical systems and suppression of nonlinearity. In: *Whither turbulence* (ed. J.L. Lumley). (Lecture Notes in Physics, vol. 357). Springer, New York, pp. 250–257
- Moin, P. (1984): Probing turbulence via large eddy simulation. *Proceedings AIAA Aerospace Sciences Meeting*
- Moin, P., Moser, R.D. (1989): Characteristic-eddy decomposition of turbulence in a channel. *J. Fluid Mech.* **200**, 471–509
- Ruelle, D., Takens, F. (1971): On the nature of turbulence. *Commun. Math. Phys.* **20**, 167–192 and **23**, 343–344
- Sell, G.R. (1989): Approximation dynamics: hyperbolic sets and inertial manifolds. (Submitted for publication)
- Sirovich, L. (1987a): Turbulence and the dynamics of coherent structures: I. *Q. Appl. Math.* **45**, 561–571
- Sirovich, L. (1987b): Turbulence and the dynamics of coherent structures: II. *Q. Appl. Math.* **45**, 561–571

- Sirovich, L. (1987c): Turbulence and the dynamics of coherent structures: III. *Q. Appl. Math.* **45**, 561–571
- Sirovich, L. (1989): Chaotic dynamics of coherent structures. *Physica D* **37**, 126–145
- Sirovich, L., Kirby, M., Winter, M. (1990): An eigenfunction approach to large scale transitional structures in jet flow. *Phys. Fluids A* **2**, 127–136
- Sirovich, L., Maxey, M., Tarman, H. (1987): Analysis of turbulent thermal convection. In: *Turbulent shear flows*, vol. 6 (eds. F. Durst et al.). Springer, New York
- Sirovich, L., Rodriguez, J.D. (1987): Coherent structures and chaos: a model problem. *Phys. Lett. A* **120**, 211–214
- Stone, E., Holmes, P. (1989): Noise induced intermittency in a model of a turbulent boundary layer. *Physica D* **37**, 20–32
- Stone, E., Holmes, P. (1990): Random perturbations of heteroclinic attractors. *SIAM J. Appl. Math.* **50**, 726–743
- Swinney, H.L., Gollub, J.P. (1981) (eds.): *Hydrodynamic instabilities and the transition to turbulence*. (*Topics in Applied Physics*, vol. 45). Springer, New York
- Temam, R. (1988): *Infinite-dimensional dynamical systems in mechanics and physics*. Springer, New York
- Tennekes, H., Lumley, J.L. (1972): *A first course in turbulence*. MIT Press, Boston, MA



# Wavelets and Applications

*Yves F. Meyer*

CEREMADE, Université Paris Dauphine, F-75775 Paris Cedex 16, France

## 1. Introduction

Some investigations conducted by (1) D. Marr in psycho-physiology of human vision, (2) J. S. Lienard in speech signal processing and (3) J. Morlet in seismic signal processing led these scientists to switch from short-time Fourier analysis to some more specific algorithms better suited to detect and analyze abrupt changes in images or signals.

These algorithms are strikingly similar and in the three of them the functions  $e^{i\omega t}$ , which have a given frequency  $\omega$ , and are the building blocks of the standard Fourier analysis, are replaced by “wavelets” which are time and frequency items and are the building blocks of “wavelets analysis”. Wavelets have a finite duration (which can be arbitrarily small) but nevertheless, should also possess a well defined average frequency.

The success of the wavelets theory is due to the remarkable formulation by A. Grossmann of J. Morlet’s ideas. Today this theory has applications in various branches of science whenever complicated interactions between events occurring at different scales appear. This happens in astrophysics [7] or in turbulence [3, 13, 14].

Independently of the above mentioned research, heavy constraints imposed by digital speech processing have led to the discovery of the so-called quadrature mirror filters. These filters also have some applications in image processing where they improve pyramidal algorithms [1].

During the fall of 1986, S. Mallat discovered that some quadrature mirror filters were the key to the construction of orthonormal wavelet bases generalizing the Haar system.

This program was completed by I. Daubechies (1987) and A. Cohen (1990) and culminated with the discovery (1989) by G. Beylkin, R. Coifman and V. Rokhlin of striking new algorithms in numerical analysis [6].

Working on submarine passive detection, J. M. Nicolas set up a new hierarchical organization of quadrature mirror filters, distinct from the one proposed by S. Mallat.

R. Coifman and the speaker proved the convergence of these schemes to new “libraries of orthonormal bases” resembling the waveforms used by J. S. Lienard.

R. Coifman and V. Wickerhauser are planning to improve speech signal compression substantially through new algorithms to select a most efficient representation from such libraries. This line of research might end the dispute between the advocates of Gabor wavelets and those of Grossmann since both types belong to the library [9].

## 2. The Windowed Fourier Transform

In signal analysis one often encounters the problem of extracting the frequency content of a signal  $f(t)$  for which one has only local information. The so-called short-time Fourier transform, or windowed Fourier transform uses cut-off functions  $g_k(t)$  which vanish when the information on  $f(t)$  is missing. The frequency content of each block  $g_k(t)f(t)$  is given by its Fourier coefficients. Let,  $g_k(t)$ , for instance, be the indicator function of the interval  $I_k = [2k\pi, 2(k+1)\pi[$  and let us expand each block  $g_k(t)f(t)$  into its Fourier series  $\sum \alpha(k, l) e^{ilt}$  on the interval  $I_k$ . It would amount to the same thing if (1), we define “trivial wavelets” by  $\sqrt{2\pi} w_{(k,l)}(t) = e^{ilt} \chi(t - 2k\pi)$ ,  $k \in \mathbb{Z}$ ,  $l \in \mathbb{Z}$ ,  $\chi(t)$  being the indicator function of  $[0, 2\pi[$ , (2) observe that these trivial wavelets form an orthonormal basis of  $L^2(\mathbb{R})$ , (3) use this basis for expanding our signal.

As everyone knows, this way of splitting into “hard blocks” produces numerical artifacts and the coefficients  $\alpha(k, l) = (2\pi)^{-1/2} \langle f, w_{k,l} \rangle$  do not give the frequency content of the signal  $f$  around  $I_k$ . To suppress these artifacts, D. Gabor decided in 1945 to replace  $\chi(t)$  by a smoother window function  $g(t)$ . Gabor wavelets are

$$(2.1) \quad w_{k,l}(t) = e^{ilt} g(t - 2k\pi), \quad l \in \mathbb{Z}, \quad k \in \mathbb{Z}$$

where  $g(t) = \pi^{-1/4} \exp(-t^2/2)$ .

But we would like the  $l^2$  norm of the wavelets coefficients to provide an energy ( $L^2$ ) estimate on the signal. This happens when  $g = \chi$ . But if  $g(t)$  satisfies the two condition  $\int_{-\infty}^{\infty} t^2 |g(t)|^2 dt < \infty$  and  $\int_{-\infty}^{\infty} \xi^2 |\hat{g}(\xi)|^2 d\xi < \infty$ , it is never the case (R. Balian, G. Battle, R.R. Coifman and S. Semmes [11]). For that reason, the definition of Gabor wavelets has been modified to

$$(2.2) \quad w_{k,l}(t) = e^{ilt} g(t - ak), \quad 0 < a < 2\pi,$$

and this over complete system is currently used in signal analysis [11].

## 3. Wavelets with Constant Shape

A function  $\psi(x)$  belonging to  $L^2(\mathbb{R}^n)$  is an *analyzing wavelet* if its Fourier transform  $\hat{\psi}(\xi)$  vanishes at 0 in a precise manner, given by the *convergence* of the following integral

$$(3.1) \quad c(\psi) = \int_{\mathbb{R}^n} |\hat{\psi}(\xi)|^2 |\xi|^{-n} d\xi.$$

If  $\int_{\mathbb{R}^n} |\psi(x)|^2 (1 + |x|)^{n+\alpha} dx < \infty$  for some  $\alpha > 0$ , then  $\psi \in L^1(\mathbb{R}^n)$  and  $c(\psi)$  is finite iff the integral of  $\psi$  vanishes. More generally  $c(\psi)$  is finite for any  $\psi$  in the Hardy space  $H^1(\mathbb{R}^n)$ .

We now pick an analyzing wavelet  $\psi$ , the mother wavelet, and decide that the other wavelets  $\psi_g$  of the family (which are used in the wavelet analysis) will only differ from  $\psi$  by their orientation, their position and their size. This can be formulated in terms of the group  $G$  with elements  $g(x) = a\varrho(x) + b$ ,  $a > 0$ ,  $b \in \mathbb{R}^n$ ,  $\varrho \in SO(n)$ . In other words,  $\psi_g(x) = a^{-n/2} \psi(g^{-1}(x))$ . The wavelet coefficients  $\alpha(g)$ ,  $g \in G$ , of  $f \in L^2(\mathbb{R}^n)$  are  $\langle f, \psi_g \rangle$  and satisfy

$$(3.2) \quad \|f\|_2^2 = (c(\psi))^{-1} \int_G |\alpha(g)|^2 dg,$$

where  $dg = a^{-n-1} da db d\varrho$  is the left-invariant Haar measure of the group  $G$ . This implies [16, 23]

$$(3.3) \quad f(x) = (c(\psi))^{-1} \int_G \alpha(g) \psi_g(x) dg.$$

Let  $\varphi$  be a radial function in the Schwartz class such that  $\hat{\varphi}(0) = 1$  and  $\hat{\varphi}(\xi) = 0$  when  $|\xi| \geq 1$ . Then J. Morlet's analyzing wavelet is

$$(3.4) \quad \psi(x) = \exp(i5x_1) \varphi(x)$$

and looks like a Gabor wavelet. But the other wavelets  $\psi$  of the Morlet family differ from Gabor wavelets by the fact that both their size and their average frequency are modified. Indeed we have

$$(3.5) \quad \psi_g(x) = a^{-n/2} e^{i5a^{-1}\nu \cdot x} \varphi\left(\frac{x-b}{a}\right)$$

where  $\nu$  is a unit vector,  $a > 0$  and  $b \in \mathbb{R}^n$ . The Fourier transform of  $\psi_g$  is contained in the ball  $|\xi - 5a^{-1}\nu| \leq a^{-1}$  which amounts to saying that the average frequency of  $\psi_g$  is  $5a^{-1}\nu$ .

On the other hand, the size of the support of  $\psi_g$  is  $O(a)$  and Morlet wavelets unlike Gabor wavelets have the sharpest spatial resolution at high frequencies, which is consistent with the Heisenberg uncertainty principle. In other words, when  $\psi$  is the Morlet analyzing wavelet, the multiscale analysis given by (3.3) also provides a multichannel analysis.

Some scientists, like D. Marr [21], have accepted a looser spectral resolution and they impose on the Fourier transform of the analyzing wavelet the following conditions:  $\hat{\psi}(\xi) \in C^\infty(\mathbb{R}^n)$  and

$$(3.6) \quad (\partial^\alpha \hat{\psi})(0) = 0 \quad \text{for } |\alpha| \leq m$$

$$(3.7) \quad \hat{\psi}(\xi) = O(|\xi|^{-N}) \quad \text{as } |\xi| \rightarrow \infty.$$

The larger  $m$  and  $N$ , the better the spectral resolution.

Scientists who apply wavelets to signal or image processing believe that the diagnostic given by the wavelet coefficients will not depend too much on the choice of  $\psi$  as long as  $\psi$  satisfies (3.6) and (3.7). This belief is supported by functional analysts who have been obtaining the same results when using either the standard Littlewood-Paley theory or one of its variants. Today it is easy to recognize that the standard Littlewood-Paley theory is a wavelet analysis based on Morlet wavelets. A. Calderón, E. Stein and G. Weiss developed a program in which a function  $f$  on  $\mathbb{R}^n$  is analyzed by the normal derivative  $\frac{\partial u}{\partial t}(x, t)$  of its harmonic extension  $u(x, t)$  to the upper half space ( $x \in \mathbb{R}^n$ ,  $t > 0$ ). This also amounts to a wavelet analysis in which the mother wavelet is  $\psi(x) = (|x|^2 - n^2)(|x|^2 + 1)^{-(n+3)/2}$ .

Grossmann's simple and elegant formalism has been directly used in astrophysics [7], in experiments on turbulence [3, 13, 14] and in many other fields of science or technology. In all these applications, the role played by wavelet analysis is to provide a better localization of small scale structures. These fine details are enhanced after being extracted from a background which is either cancelled or strongly attenuated by a correctly tuned wavelet.

Discrete versions of (3.3) would be  $f(x) = \sum_{\lambda \in \Lambda} c(\lambda)\psi_\lambda(x)$  where  $\Lambda$  is a suitable discrete subset of  $G$ . J. Morlet proposed  $\lambda(x) = 2^{-\alpha j}(r(x) + \beta k)$ ,  $j \in \mathbb{Z}$ ,  $k \in \mathbb{Z}^n$ ,  $r \in F$ , where  $\alpha > 0$ ,  $\beta > 0$  are small enough and  $F$  is a finite set of rotations. I. Daubechies proved in [11] that

$$\left\| \sum_{\lambda \in \Lambda} c(\lambda)\psi_\lambda(x) \right\|_{L^2(\mathbb{R}^n)} \leq C \left( \sum_{\lambda \in \Lambda} |c_\lambda|^2 \right)^{1/2}$$

and that the corresponding mapping from  $l^2(\Lambda)$  into  $L^2(\mathbb{R}^n)$  is onto when  $\alpha$  and  $\beta$  are small enough. The values of  $\alpha$  and  $\beta$  depend strongly on the choice of  $\psi$  which should satisfy (3.6) and (3.7) [22].

Daubechies' theorem implies that a canonical decomposition  $f(x) = \sum_{\lambda \in \Lambda} c(\lambda)\psi_\lambda(x)$  exists in which the coefficients are given by  $c(\lambda) = \langle f, \tilde{\psi}_\lambda \rangle$ . Unfortunately these "dual wavelets"  $\tilde{\psi}_\lambda$  might be badly behaved in terms of size and regularity. When this happens, it prevents us from using (in G. Weiss terminology) the "atomic decomposition"  $f(x) = \sum c(\lambda)\psi_\lambda(x)$  in situations other than the "trivial  $L^2$  setting". But if one is limited to the  $L^2$  theory, the Haar system (1909) will provide the most efficient wavelet analysis, since it uses the simplest analyzing wavelet and since it is an orthonormal basis.

All these difficulties would disappear if  $\psi_\lambda$ ,  $\lambda \in \Lambda$ , forms an orthonormal basis of  $L^2(\mathbb{R}^n)$ .

#### 4. Pyramidal Algorithms, Quadrature Mirror Filters and Orthonormal Wavelets

We fix an integer  $N \geq 1$  and assume we are given two trigonometric polynomials  $m_0(\theta) = h_0 + h_1 e^{i\theta} + \cdots + h_{2N-1} e^{i(2N-1)\theta}$ ,  $m_1(\theta) = g_0 + g_1 e^{i\theta} + \cdots + g_{2N-1} e^{i(2N-1)\theta}$  such that  $m_0(0) = 1$ ,  $|m_0(\theta)|^2 + |m_0(\theta + \pi)|^2 = 1$  and  $m_1(\theta) = e^{i(2N-1)\theta} \overline{m_0(\theta + \pi)}$ .

We fix these polynomials and for every  $\delta > 0$  we consider the corresponding operators  $F_0 : l^2(\delta\mathbb{Z}) \rightarrow l^2(2\delta\mathbb{Z})$ ,  $F_1 : l^2(\delta\mathbb{Z}) \rightarrow l^2(2\delta\mathbb{Z})$  which are defined by

$$(4.1) \quad F_0\{x_{(k\delta)}\}(2l\delta) = \sqrt{2} \sum_{-\infty}^{\infty} h_k x_{(k+2l)\delta},$$

$$(4.2) \quad F_1\{x_{(k\delta)}\}(2l\delta) = \sqrt{2} \sum_{-\infty}^{\infty} g_k x_{(k+2l)\delta}.$$

These two operators are called “quadrature mirror filters” and have the following remarkable properties. The operator  $F = (F_0, F_1) : l^2(\delta\mathbb{Z}) \rightarrow l^2(2\delta\mathbb{Z}) \times l^2(2\delta\mathbb{Z})$  is a unitary isomorphism and both adjoints  $F_0^* : l^2(2\delta\mathbb{Z}) \rightarrow l^2(\delta\mathbb{Z})$ ,  $F_1^* : l^2(2\delta\mathbb{Z}) \rightarrow l^2(\delta\mathbb{Z})$  are partial isometries. The ranges of  $F_0^*$  and  $F_1^*$  are orthogonal in  $l^2(\delta\mathbb{Z})$  and, finally, one has

$$(4.3) \quad I = F_0^*F_0 + F_1^*F_1.$$

We now consider the increasing sequence  $\Gamma_j = 2^{-j}\mathbb{Z}$  of lattices ( $0 \leq j$ ) together with the corresponding partial isometries  $F_0^* : l^2(\Gamma_j) \rightarrow l^2(\Gamma_{j+1})$ . It makes sense to compose these operators and we are led to study the asymptotic behavior of  $(F_0^*)^j : l^2(\Gamma_0) \rightarrow l^2(\Gamma_j)$ . By construction, this operator commutes with integral translations  $\tau_k$ ,  $k \in \mathbb{Z}$ . Let  $\varepsilon_k \in l^2(\Gamma_0)$  be defined by  $\varepsilon_k(l) = 0$  if  $l \neq k$ , 1 if  $l = k$ . Then  $(F_0^*)^j(\varepsilon_k)$ ,  $k \in \mathbb{Z}$ , is an orthonormal sequence in  $l^2(\Gamma_j)$  and we would like to know if, in some sense, this sequence converges to an orthonormal sequence of the form  $\varphi(x - k)$ ,  $k \in \mathbb{Z}$ , where  $\varphi$  belongs to  $L^2(\mathbb{R})$ . The convergence procedure is defined in the following way. To each element  $f_j(k2^{-j})$  in  $l^2(\Gamma_j)$  we attach the corresponding atomic measure  $\sigma_j = 2^{-j/2} \sum f_j(k2^{-j})\delta_{(k2^{-j})}$  where  $\delta_a$  is the Dirac mass at the point  $a$ . We next define the convergence of the sequence  $f_j$  by the weak convergence of the corresponding sequence  $\sigma_j$ .

**Theorem 2** (A. Cohen, 1989). *The two following conditions are equivalent*

- (4.4) *the discrete orthonormal sequences  $(F_0^*)^j(\varepsilon_k)$ ,  $k \in \mathbb{Z}$ , converge to  $\varphi(x - k)$ ,  $k \in \mathbb{Z}$  and this sequence is orthonormal in  $L^2(\mathbb{R})$ .*
- (4.5) *for  $0 < \theta < 2\pi$ ,  $\lim_{j \rightarrow \infty} m_0(\theta)m_0(2\theta)\dots m_0(2^j\theta) = 0$ .*

If (4.5) is satisfied, the function  $\varphi \in L^2(\mathbb{R})$  which is called the *scaling function* can be characterized by an other property. Consider the functional equation

$$(4.6) \quad \varphi(x) = 2 \sum_0^{2N-1} h_k \varphi(2x - k),$$

where we impose the condition  $\varphi \in L^1(\mathbb{R})$  and  $\int \varphi(x) dx = 1$ . Then (4.6) admits a unique solution which is precisely the function  $\varphi$  of (4.4). Moreover the Fourier transform  $\hat{\varphi}(\xi)$  is given by

$$(4.7) \quad \hat{\varphi}(\xi) = m_0(\xi/2)\dots m_0(\xi/2^j)\dots$$

We next define  $\psi \in L^2(\mathbb{R})$  by  $\psi(x) = 2 \sum_0^{2N-1} g_k \varphi(2x - k)$  or equivalently by

$$(4.8) \quad \hat{\psi}(\xi) = m_1(\xi/2) m_0(\xi/4) \dots m_0(\xi/2^j) \dots$$

If (4.5) is satisfied, then  $2^{j/2}\psi(2^j x - k)$ ,  $j \in \mathbb{Z}$ ,  $k \in \mathbb{Z}$ , is an orthonormal basis of  $L^2(\mathbb{R})$ .

It remains, for a given  $N$ , to choose  $m_0(\theta)$  and  $m_1(\theta)$  carefully in order to obtain a smooth wavelet  $\psi(x)$ . One way to proceed is the following. There exists (F. Riesz)  $m_0(\theta) = h_0 + \dots + h_{2n-1} e^{i(2N-1)\theta}$  such that  $m_0(0) = 1$  and  $|m_0(\theta)|^2 = c_N \int_\theta^\pi (\sin x)^{2N-1} dx$  ( $c_N > 0$  will permit  $m_0(0) = 1$ ). We then have [12].

**Theorem 3** (I. Daubechies, 1987). *There exists a constant  $\alpha > 0$  such that for each  $N \geq 2$ , the corresponding functions  $\varphi$  and  $\psi$  will belong to  $C^{\alpha N}$ .*

This construction also gives orthonormal wavelet bases in several dimensions. For the sake of simplicity, we stick to  $n = 2$ . We then consider the three wavelets  $\psi_1(x, y) = \psi(x) \varphi(y)$ ,  $\psi_2(x, y) = \varphi(x) \psi(y)$  and  $\psi_3(x, y) = \psi(x) \psi(y)$ . The full collection  $2^j \psi_q(2^j x - k, 2^j y - l)$ ,  $q = 1, 2$  or  $3$ ,  $j \in \mathbb{Z}$ ,  $k \in \mathbb{Z}$ ,  $l \in \mathbb{Z}$  is an orthonormal basis of  $L^2(\mathbb{R}^2)$ .

When compared to Fourier series expansions, orthonormal wavelets expansions provide a much deeper insight into the local or global properties of the function to be analyzed [17, 18].

For example, J. O. Strömberg proved in 1980 that orthonormal wavelets form an unconditional basis for the Hardy space  $H^1(\mathbb{R}^n)$ . A suitable regrouping of the wavelet expansion of a function  $f$  in  $H^1(\mathbb{R}^n)$  yields its atomic decomposition [22].

## 5. Wavelets Packets

Wavelet analysis gives its best performance when it is applied to signals with abrupt changes or to functions which have simple discontinuities on smooth surfaces and which are smooth elsewhere. On the other hand, wavelet analysis gives its worst performance on stationary signals.

The speech signal obviously contains these two components and one would like to switch freely from wavelet analysis to windowed Fourier analysis in speech signal processing. R. Coifman and his co-workers have constructed a library of orthonormal bases in which can be found I. Daubechies orthonormal wavelets, a second orthonormal basis resembling the Gabor wavelets as well as many other bases.

Let us fix some notations to describe this library. We start with two quadrature mirror filters  $F_0 : l^2(\mathbb{Z}) \rightarrow l^2(2\mathbb{Z})$  and  $F_1 : l^2(\mathbb{Z}) \rightarrow l^2(2\mathbb{Z})$  as in Theorem 3. Let  $\mathcal{I}$  be the collection of all dyadic intervals  $I = [l2^j, (l+1)2^j]$ ,  $l \in \mathbb{N}$ ,  $j \in \mathbb{Z}$ , which are contained in  $[0, \infty)$ . To each  $I \in \mathcal{I}$  we attach a closed subspace  $W_I \subset L^2(\mathbb{R})$  and a function  $w_I(x)$  which are uniquely defined by the four following properties

$$(5.1) \quad w_I(x - k2^{-j}), \quad k \in \mathbb{Z}, \text{ is an orthonormal basis of } W_I;$$

- (5.2) if  $I_0$  is the left half of  $I$  and  $I_1$  the right half, then  $W_I$  is the direct orthogonal sum of  $W_{I_0}$  and  $W_{I_1}$ ;
- (5.3) if  $f(x) = \sum_{-\infty}^{\infty} \alpha_I(k) w_I(x - k2^{-j})$   
 $= \sum_{-\infty}^{\infty} \beta_{I_0}(2k) w_{I_0}(x - 2k2^{-j}) + \sum_{-\infty}^{\infty} \gamma_{I_1}(2k) w_{I_1}(x - 2k2^{-j}),$   
then  $\beta_{I_0} = F_0(\alpha_I)$  and  $\gamma_{I_1} = F_1(\alpha_I)$ ;
- (5.4) if  $I = [0, 2^j]$ , then  $W_I = V_j$  and  $w_I(x) = 2^{j/2} \varphi(2^j x)$ .

It is easy to compute  $w_I(x)$ . We write  $I = [l2^j, (l+1)2^j)$ ,  $l = \varepsilon_0 + \varepsilon_1 2 + \cdots + \varepsilon_q 2^q + \dots$  where  $\varepsilon_q = 0$  or 1 and we have

$$(5.5) \quad \hat{w}_I(\xi) = 2^{-j/2} m_{\varepsilon_0}(\xi/2^{j+1}) m_{\varepsilon_1}(\xi/2^{j+2}) \dots$$

With these notations, the library of orthonormal bases generated by the quadrature mirror filters  $(F_0, F_1)$  is described by the following theorem.

**Theorem 4.** Let  $\mathcal{J}_* \subset \mathcal{J}$  be any collection of dyadic intervals  $I \subset [0, \infty)$  with the property that, excepting a set  $D$  which is either finite or denumerable, each  $x \in [0, \infty[, x \notin D$ , belongs to one (and only one) interval  $I \in \mathcal{J}_*$ .

Then the corresponding family  $w_I(x - k2^{-j})$ ,  $k \in \mathbb{Z}$ ,  $I = [l2^j, (l+1)2^j) \in \mathcal{J}_*$ , is an orthonormal basis of  $L^2(\mathbb{R})$ .

When  $\mathcal{J}_*$  is the obvious collection of the dyadic intervals  $[2^j, 2^{j+1})$ ,  $j \in \mathbb{Z}$ , this basis happens to be the one described in theorem 3 and when  $\mathcal{J}_*$  is the collection of  $[l, l+1)$ ,  $l \in \mathbb{N}$ , the corresponding  $w_I(x)$  will be denoted by  $w_l(x)$  and resemble Gabor wavelets. Moreover, if  $m_0(\xi) = \frac{1+e^{i\xi}}{2}$  and  $m_1(\xi) = \frac{1-e^{i\xi}}{2}$ , the orthonormal basis  $w_l(x - k)$ ,  $l \in \mathbb{N}$ ,  $k \in \mathbb{Z}$ , is the well known Walsh system which is widely used in signal processing. When  $N \geq 2$  and  $m_0(\xi)$  is chosen following I. Daubechies, the corresponding orthonormal basis  $w_l(x - k)$ ,  $l \in \mathbb{N}$ ,  $k \in \mathbb{Z}$ , should be compared to a smooth version of the Walsh system.

In their work in speech signal compression, R. R. Coifman and V. Wickerhauser are using this full library of bases together with an entropy criterion for selecting the specific basis among the library which provides the “best” expansion [9].

## References

1. E.H. Adelson, R. Hingorani, E. Simoncelli: Orthogonal pyramid transforms for image coding. SPIE Visual Communications and Image Processing II. **845** (1987)
2. M. Antonini, M. Barlaud, I. Daubechies, P. Mathieu: Image coding using wavelet transform. Lassy, Nice-Sophia Antipolis, 06560 Valbonne, France
3. F. Argoul, A. Arnéodo, Y. Gagne, G. Grasseau, U. Frisch, E.J. Hopfinger: Wavelet analysis of turbulence reveals the multifractal nature of the Richardson cascade. Nature **338**, no. 6210 (1989)
4. G. Battle: Heisenberg proof of the Balian-Low theorem. Math. Phys. Letters (1990)
5. G. Battle, P. Federbush: Ondelettes et phase cluster expansions: a vindication. Comm. Math. Phys. **109** (1987) 417–419

6. G. Beylkin, R. Coifman, V. Rokhlin, Fast Wavelet Transforms and Numerical Algorithms, I. Research Report YALEU/DCS/RR-696, 1989
7. A. Bijaoui, G. Mars, E. Slezak: Identification of structures from galaxy counts: use of the wavelet transform. *Astron. Astrophys.* **227** (1990) 301–316
8. A. Cohen: Ph.D. Dissertation, CEREMADE, Université Paris-Dauphine, 75775 Paris, Cedex 16
9. R. Coifman, V. Wickerhauser: Best-adapted wave packet bases. Numerical Algorithms Research Group, Dpt. Mathematics, Yale, New Haven, CT 06520, USA
10. J.M. Combes, A. Grossmann Ph. Tchamitchian: Time frequency methods and phase space. Proceedings of the International conference, Marseille, France, Dec. 14–18, 1987
11. I. Daubechies: The wavelet transform, time-frequency localization and signal analysis. *IEEE Trans. Inf. Theory* (1990)
12. I. Daubechies: Orthogonal bases of compactly supported wavelets. *Comm. Pure Appl. Math.* **41** (1988) 909–996
13. M. Farge, M. Holschneider, J.F. Colonna: Wavelet analysis of coherent structures in two-dimensional turbulence. *Topological Fluid Mechanics*, ed. K. Moffat. Cambridge University Press, 1989, pp. 765–776
14. M. Farge, G. Rabreau: Wavelet transform to detect and analyze coherent structures in two-dimensional turbulent flows. *C. R. Acad. Sci. Paris* **307** Ser. II (1988) 1479–1486
15. D. Gabor: Theory of communication. *J. Inst. Elec. Eng. (London)* **93** III (1946) 429–457
16. A. Grossmann, J. Morlet: Decomposition of Hardy functions into square integrable wavelets of constant shape. *SIAM J. Math. Anal.* **15** (1984) 723–736
17. M. Holschneider, Ph. Tchamitchian: Pointwise Analysis of Riemann's “nondifferentiable” function. *Invent. math.* (to appear)
18. S. Jaffard: Exposants de Hölder en des points donnés et coefficients d'ondelettes. *C. R. Acad. Sci. Paris* **308** Ser. I (1989) 79–81
19. P.G. Lemarié: Les ondelettes en 1989. (*Lecture Notes in Mathematics*, vol. 1438). Springer, Berlin Heidelberg New York 1990
20. S. Mallat: Review of multifrequency channel decompositions of images and wavelet models. *IEEE Trans. Acoustics, Speech and Signal Processing* **37**, no. 12 (1989)
21. D. Marr: Vision: A computational investigation into the human representation and processing of visual information. Freeman, New York 1982
22. Y. Meyer: Ondelettes et opérateurs, tomes I, II et III. Hermann, Paris 1990
23. R. Murenzi: Wavelet transforms associated to the  $n$ -dimensional euclidean group with dilations. Preprint Institut de Physique Théorique, Université Catholique de Louvain, B-1348 Louvain-la-Neuve, Belgium

# Pattern Formation in Reaction-Diffusion Systems

Masayasu Mimura

Department of Mathematics, Faculty of Science, Hiroshima University  
Hiroshima, 730 Japan

## 1. Introduction

Reaction-diffusion approach has been used to explain pattern formation arising in neurobiology, chemical physics, population ecology, developmental biology and other fields. Despite its simple structure, a class of reaction-diffusion systems exhibit a lot of spatial and spatio-temporal patterns. Some of these patterns in a reacting and diffusing medium can be often observed by internal layers or interfaces which are boundaries between qualitatively different states in the system. Such interfaces exhibit a variety of geometrical patterns such as rotating patterns in the Belousov-Zhabotinsky reagent [Wi], dendritic patterns in solidifications [Ca], pigmentation patterns on shells [MK] and animal coat marking [Mu], for instance.

The term “reaction-diffusion equations” is usually taken to mean the following semilinear parabolic equations:

$$\frac{\partial u}{\partial t} = D \Delta u + F(u), \quad (1.1)$$

where  $u(t, x) = (u_1, u_2, \dots, u_n)(t, x)$  means the concentration, density and other physical component with time and space variables  $t$  and  $x$ ,  $D$  is an  $n \times n$  nonnegative matrix and in most cases, it is a diagonal one and  $F$  is the reaction term.

Among so many reaction-diffusion equation models, we are concerned with activator-inhibitor systems which arise in modelling of morphogenesis ([GM]). The most simple and suggestive system is the following two-component model:

$$\begin{cases} \frac{\partial u}{\partial t} = d_1 \Delta u + f(u, v) \\ \frac{\partial v}{\partial t} = d_2 \Delta v + \delta g(u, v), \end{cases} \quad (1.2)$$

where  $u$  and  $v$  are called respectively the activator and its inhibitor in morphogenesis [GM] or the propagator and its controller in excitable media [Fi1]. Here  $d_1$  and  $d_2$  are the diffusion rates of  $u$  and  $v$ ,  $\delta$  is the ratio of reaction rates. The nonlinearities of  $f$  and  $g$  in which we are interested in this paper are restricted to two types in

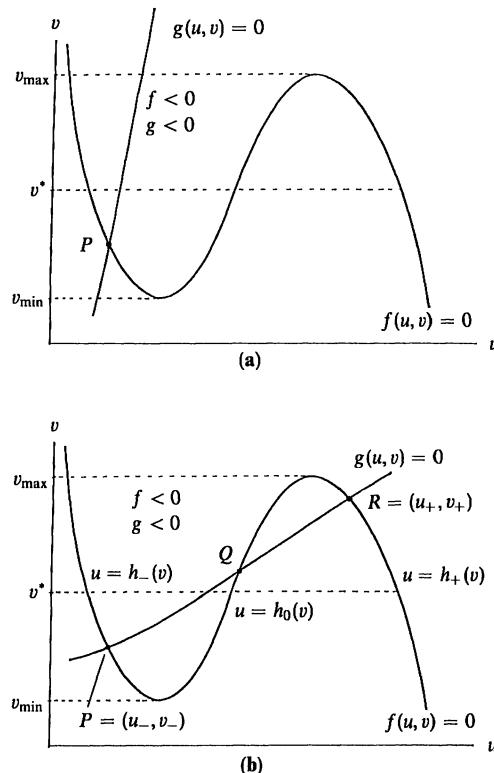


Fig. 1

Figs. 1a and b. A prototype of  $f$  and  $g$  is

$$\begin{cases} f(u, v) = u(1-u)(u-a) - v \\ g(u, v) = u - \gamma v \end{cases} \quad (1.3)$$

with constants  $0 < a < 1$  and  $\gamma > 0$ . When  $v$  is totally absent in (1.2) with (1.3), it is called the time dependent Ginzburg-Landau equation and when  $v$  does not diffuse, it reduces to the well-known FitzHugh-Nagumo equation which models a nerve impulse propagating along the axon([NAY]).

The kinetic system of (1.2) is given by

$$\begin{cases} \frac{du}{dt} = f(u, v) \\ \frac{dv}{dt} = \delta g(u, v), \end{cases} \quad (1.4)$$

which is called the Bonhoeffer-Van der Pol equation when  $f$  and  $g$  take (1.3).

For the case in Fig. 1a, (1.4) has the only one equilibrium state  $P$ , which is globally stable. In this case, a small disturbance from this state is rapidly damped, while a large disturbance wanders far from the state but then eventually returns to the state. Because of this feature, the state  $P$  is called a rest state in excitable media. On the other hand, for the case in Fig. 1b, (1.4) has three equilibria  $P$ ,  $Q$  and  $R$  where  $P$  and  $R$  are stable and  $Q$  is unstable. That is, (1.4) is called a monostable system for the former case, while it is a bistable one for the latter case.

We consider the situation where diffusion terms are present in (1.2) and one component  $u$  diffuses slower than the other  $v$ , moreover,  $u$  reacts much faster than  $v$ . More specifically, we introduce the new parameters  $\varepsilon$ ,  $\tau$  and  $D$  instead of  $d_1$ ,  $d_2$  and  $\delta$  through

$$\varepsilon = \sqrt{d_1}, \quad \tau = \delta/\sqrt{d_1} \quad \text{and} \quad D = d_2/\delta.$$

Moreover, we specify the nonlinearities of  $f$  and  $g$  as (1.8) for simplicity only, though the results which will be stated later are valid for more general nonlinearities under appropriate conditions. We thus rewrite (1.2) as

$$\begin{cases} \varepsilon\tau \frac{\partial u}{\partial t} = \varepsilon^2 \Delta u + f(u) - v \\ \frac{\partial v}{\partial t} = D\Delta v + u - \gamma v \end{cases} \quad (1.5)$$

with  $f(u) = u(1-u)(u-a)$ , where  $\varepsilon$  is sufficiently small and  $\tau$  and  $D$  are of the order  $O(1)$  compared with  $\varepsilon$ .

Assuming that  $\varepsilon$  is sufficiently small, we use singular perturbation techniques to study the existence and stability of stationary pulse solutions in a monostable system and traveling front solutions in a bistable one in Sections 2 and 3, respectively. Both solutions exhibit internal layers with width  $O(\varepsilon)$ . In Section 4, in order to study the dynamics of such layers, we derive the equation of motion for interfaces in the limit  $\varepsilon \downarrow 0$ . Finally, we would like to give some remarks on our system in Section 5.

*Acknowledgement.* I have benefitted from the discussions with my colleagues, H. Fujii, Y. Nishiura, H. Ikeda, T. Tsujikawa, R. Kobayashi and T. Ohta.

## 2. Localized Patterns in a Monostable Medium

In this section, we consider (1.5) under the situation where there is only one trivial (constant) equilibrium state  $P = (0, 0)$  as in Fig. 1a. In addition, we assume

$$\int_0^1 f(u) du > 0$$

under which the equilibrium state  $P$  is asymptotically stable. This situation is realizable when  $a$  and  $\gamma$  are chosen to satisfy  $0 < a < 1/2$  and  $0 < \gamma < \max(u-a)(1-u)$ .

We treat the system (1.5) in the whole domain  $R^N$ . The boundary condition at infinity is

$$\lim_{|x| \rightarrow \infty} (u, v) = (0, 0). \quad (2.1)$$

For this problem, there arises naturally a question of interest: Are there any nontrivial equilibrium states? Intuitively we can imagine the following situation: Suppose that there is a local disturbance in  $u$ . If it is not so small, then it possibly forms into a large peak and expands, as an activator, but its expanding may be stopped due to the relatively faster diffusion of the inhibitor. This argument suggests the possibility of the localization of two components, if there exists a suitable balance between them.

Motivated by this suggestion, we are interested in the existence problem of (hopefully stable) nonconstant equilibrium states of (1.5) under the boundary condition (2.1).

For one dimensional case, singular perturbation techniques [Fi2] or shooting arguments [EHT] can be applied to (1.5), (1.6) if  $\varepsilon$  is sufficiently small.

**Theorem 1.** *There is  $\varepsilon_0$  such that the stationary problem of (1.5), (2.1) has two different types of 1D-pulse solutions  $(\bar{u}_\varepsilon, \bar{v}_\varepsilon)$  and  $(\underline{u}_\varepsilon, \underline{v}_\varepsilon)$ , which are symmetric with  $x = 0$ , for  $0 < \varepsilon < \varepsilon_0$  (the profiles of these solutions are Fig. 2).*

We now address the following questions: (1) Is there any other pulse-solution for small  $\varepsilon$ ? (2) Is there any pulse solution for not small  $\varepsilon$ ? (3) How is the stability of the equilibrium solutions  $(\bar{u}_\varepsilon, \bar{v}_\varepsilon)$  and  $(\underline{u}_\varepsilon, \underline{v}_\varepsilon)$ ? Unfortunately, the first two equations (1) and (2) have not yet been answered except for the special case when  $f$  is piecewise-linear

$$f(u, v) = -1 + H(u - a) - v$$

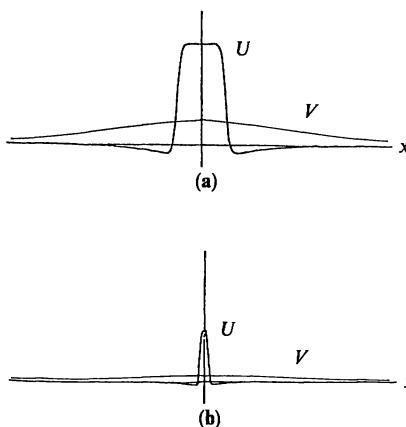


Fig. 2

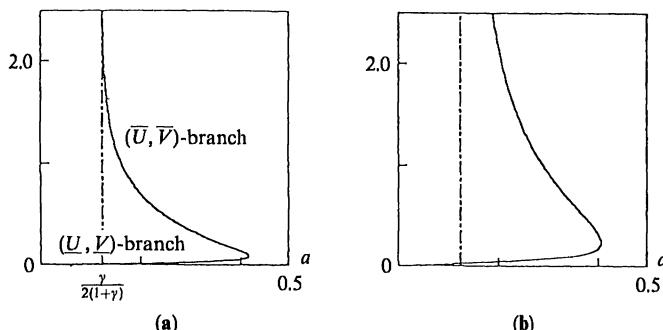


Fig. 3

with constants  $a$  and  $\gamma$  satisfying  $0 < a < 1/2$  and  $0 < \gamma < 2a/(1 - 2a)$ , where  $H(x)$  is the Heaviside step function. For this case, the global structure of symmetric pulse solutions can be completely understood ([OMK]). For instance, Figs. 3a and b show respectively 1D- and 3D-radially symmetric solutions. That is, for suitably fixed  $a$  and  $\gamma$ , there is the critical value  $\varepsilon_c$  such that there are exactly two different solutions for  $0 < \varepsilon < \varepsilon_c$ , where the upper and lower branches correspond to  $(\bar{u}_\varepsilon, \bar{v}_\varepsilon)$  and  $(\underline{u}_\varepsilon, \underline{v}_\varepsilon)$ , respectively, while there are no solutions for  $\varepsilon_c < \varepsilon$ .

For the third question (3), we have recently obtained the following:

**Theorem 2 [NM].** *Let  $(\bar{u}_\varepsilon, \bar{v}_\varepsilon)$  and  $(\underline{u}_\varepsilon, \underline{v}_\varepsilon)$  be the equilibrium solutions of (1.5), (2.1) which are given in Theorem 1. There is the critical value  $\tau_c$  such that  $(\bar{u}_\varepsilon, \bar{v}_\varepsilon)$  is asymptotically stable (except for translation free) for  $\tau_c < \tau$ , while it is unstable through Hopf bifurcation for  $0 < \tau < \tau_c$ . On the other hand,  $(\underline{u}_\varepsilon, \underline{v}_\varepsilon)$  is unstable for any  $\tau > 0$ .*

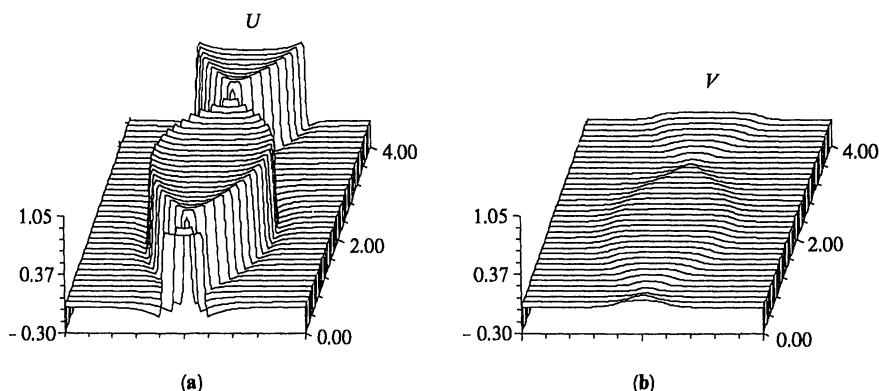


Fig. 4

**Remark [[NM]].** When  $\tau$  decreases, the destabilization of  $(\bar{u}_\epsilon, \bar{v}_\epsilon)$  occurs and there appears a bifurcating solution which exhibits oscillating internal layers as in Fig. 4.

It is quite interesting to consider the stability of pulse solutions in higher dimensional space. For the piecewise nonlinearity, it has been recently shown in [OMK] that even if  $\tau$  is fixed to satisfy  $\tau_c < \tau$ , the stability of 2D- and 3D-solutions is different from the 1D-one. For higher dimensional cases, the stability crucially depends on parameters  $a$  and  $\gamma$ . For instance, when  $\gamma$  is suitably fixed, there is the critical value  $a_c$  such that the solution is stable for  $a_c < a$ , while it is unstable through radially symmetric breaking. For more general nonlinearities, numerical simulations show that such destabilization also occurs ([OMK]). Its rigorous treatment is in progress.

### 3. Traveling Fronts in a Bistable Medium

In this section, we consider (1.5) in the case when there are two equilibrium states which are stable as in Fig. 1b. We simply write them as  $P : (u_-, v_-)$  and  $R : (u_+, v_+)$ . This situation is realizable when  $0 < a < 1/2$  and

$$\gamma > \gamma_0 = q / \{(2 - a - \sqrt{a^2 - a + 1})(1 - 2a + \sqrt{a^2 - a + 1})\}$$

Under this situation, we consider one dimensional traveling front solutions which correspond to propagating transition from one state  $P$  to the other  $R$ . This kind of problem occurs in the study of diffusive waves in population genetics, combustion theory, chemical reaction and population dynamics.

These solutions can be represented by  $(u, v)(z)$  with  $z = x + ct$  where  $c$  is the velocity. Thus, (1.5) is written as

$$\begin{cases} \epsilon^2 u_{zz} - \epsilon c \tau u_z + f(u) - v = 0 \\ v_{zz} - cv_z + u - \gamma v = 0 \end{cases} \quad z \in R, \quad (3.1)$$

where we may take  $D = 1$  without loss of generality. The corresponding boundary conditions are

$$(u, v)(\pm\infty) = (u_\pm, v_\pm). \quad (3.2)$$

Our problem is to find  $c$  such that (3.1) has a heteroclinic orbit  $u(z)$  connecting  $(u_+, v_+)$  and  $(u_-, v_-)$  at the infinities  $z = \pm\infty$ .

For a special case when  $\gamma$  tends to  $\infty$ , (3.1), (3.2) formally reduces to the scalar problem for  $u$

$$\begin{cases} \epsilon^2 u_{zz} - \epsilon c \tau u_z + f(u) = 0, \\ u(\pm\infty) = u_\pm \end{cases} \quad z \in R \quad (3.3)$$

because  $v \equiv 0$ . For this problem, it is already known that there is only one traveling front solution  $u(z)$  with the unique velocity  $c$  and it is stable (except for translation free) ([FM], for instance). We easily find that when  $\epsilon$  is sufficiently small, this solution exhibits a single internal layer with width  $O(\epsilon)$ .

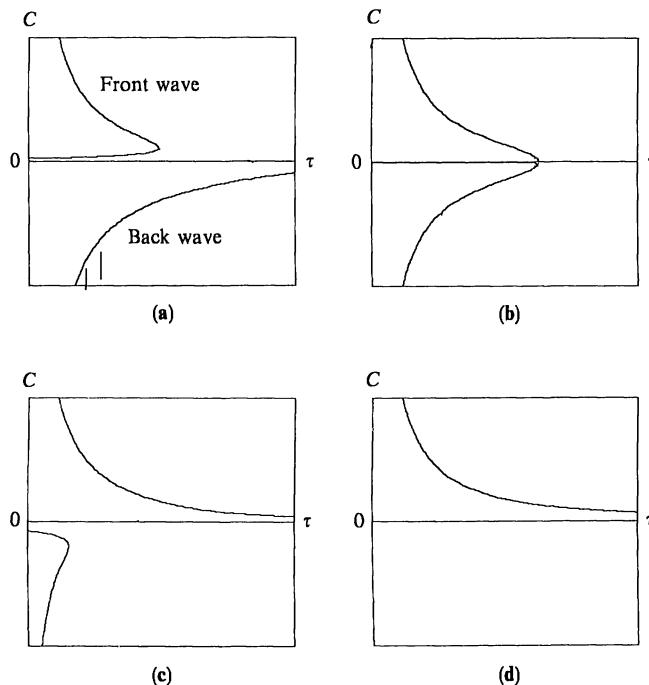
We consider the problem as to how many traveling fronts exist and how their stability are for (3.1), (3.2). Before stating the result, we introduce the following two critical values in addition to  $\gamma_0$ :

$$\gamma_1 = 9/\{(2-a)(1-2a)\}$$

and

$$\gamma_2 = 9/(1+a+\sqrt{3(a^2-a+1)})/\{(1+a)(2-a)(1-2a)\}$$

When  $\varepsilon$  is sufficiently small, singular perturbation techniques can apply to (3.1), (3.2) so that the global bifurcation pictures with respect to  $\tau$  are drawn for some values of  $\gamma$ , as in Fig. 5 ([IMN]). We find that the number of traveling front solutions depend upon the parameters  $\tau$  and  $\gamma$ , which is totally different from the scalar version. Fig. 5a is the structure for  $\gamma_0 < \gamma < \gamma_1$ , Fig. 5b is for  $\gamma = \gamma_1$ , Fig. 5c is for  $\gamma_1 < \gamma < \gamma_2$  and Fig. 5d is for  $\gamma_2 < \gamma$ . For the special case when  $\gamma = \gamma_1$  so that the kinetics possess odd symmetry, there exist an odd symmetric standing front solution (which has zero velocity) for any  $\tau$  and a pitchfork bifurcation occurs at  $\tau = \tau_c$  such that the trivial solution is stable for  $\tau_c < \tau$ , while it is unstable for  $0 < \tau < \tau_c$ . When  $\gamma$  is slightly different from  $\gamma_1$ , this symmetric structure is deformed into the imperfection structures in which one limit point appears.



**Fig. 5**

We call the front with positive velocity a front wave, while the one with negative velocity a back wave. For the stability of these solutions in one dimensional case, we have

**Theorem 3 (NMIF]).** Consider the solution branches in Fig. 5.

- (1) (Fig. 5a) The fast front wave and the back wave are stable, while the slow front wave is unstable;
- (2) (Fig. 5b) The front and back waves are both stable;
- (3) (Fig. 5c) The front wave and the fast back wave are stable, while the slow back wave is unstable;
- (4) (Fig. 5d) The front wave is stable.

The fourth case corresponds to the scalar problem which was already noted.

A natural question arises whether these planar traveling front solutions are stable or not in higher dimensional spaces? Recently we have found that the stability crucially depends on the value of  $\tau$ . The result will be stated in a forthcoming paper.

#### 4. Interfacial Dynamics

We have shown the stability as well as existence of stationary pulse solutions and traveling front solutions which possess internal layers with width  $0(\varepsilon)$ . From an application view point, we address the following question: When  $\varepsilon$  is sufficiently small, how does a solution of the initial-boundary value problem of (1.5), (2.1) or (1.5), (3.2) evolve into these patterns? In order to answer it, we derive the approximating evolution equation from the original reaction-diffusion system (1.5) when  $\varepsilon$  tends to zero. In this case, one could expect that an internal layer becomes an interface in the limit  $\varepsilon \downarrow 0$ . First, we consider a bistable scalar equation of (1.5) when  $v$  is assumed to be constant, say  $v = q$  such that  $f(u) = q$  has three zeros  $h_-(q)$ ,  $h_0(q)$  and  $h_+(q)$ .

$$\varepsilon\tau \frac{\partial u}{\partial t} = \varepsilon^2 \Delta u + f(u) - q. \quad (4.1)$$

We consider (4.1) in the whole plane  $R^2$  for simplicity. When  $\varepsilon$  is sufficiently small, one could expect that the process consists of two different stages. The first stage is that even if the initial data is smooth, the solution  $u(t, x)$  tends, in a short time, to one stable equilibrium state  $u = h_+(q)$  in a region where  $u(0, x) > h_0(q)$  or to the other stable state  $u = h_-(q)$  in a region where  $u(0, x) < h_0(q)$ , since the system is a bistable one. This indicates that there appear internal layers with width  $0(\varepsilon)$  which separate the plane into two different subplanes ([FH]). In the limit  $\varepsilon \downarrow 0$ , it implies the appearance of interfaces in the plane. The next stage is that such interfaces propagate. Suppose that the interface is described by a gentle curve  $\Gamma(t)$  in a way that  $R^2 \setminus \Gamma(t) = \Omega_+(t) \cup \Omega_-(t)$  where the relations  $u = h_{\pm}(q)$  hold in the regions  $\Omega_{\pm}(t)$ . The motion of  $\Gamma(t)$  is described by the following evolution equation ([KT], for instance):

$$\tau V = \{-c(q) + \varepsilon K\}n \quad (4.2)$$

where  $V$  is the normal velocity of the interface,  $K$  is the mean curvature at  $\Gamma(t)$ ,  $n$  is the unit normal vector of  $\Gamma(t)$  pointing from  $\Omega_+(t)$  to  $\Omega_-(t)$  and  $c(q)$  is the velocity of the 1D-traveling front solution of the bistable equation

$$\begin{cases} \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + f(u) - q, & t > 0, \quad x \in R \\ \lim_{x \rightarrow \pm\infty} u = h_{\pm}(q) \end{cases}$$

It is known that  $c(q)$  is explicitly represented as

$$c(q) = \frac{h_+(q) - 2h_0(q) + h_-(q)}{\sqrt{2(h_+(q) - h_-(q))}}$$

Recently, for the study of the curvature effect on the motion of interfaces, (4.2) has been investigated in the mathematical community ([CGG], [ES], [Gr] and the references therein).

Come back to our system of equations (1.5) in  $R^2$ . Keeping the scalar equation (4.1) in mind, one could expect that the process also consists of two stages: The first stage is the appearance of internal layers or interfaces in the limit  $\varepsilon \downarrow 0$ . In the region where is away from interfaces, we may put  $\varepsilon = 0$  in (1.5) so that it becomes

$$\frac{\partial v}{\partial t} = D\Delta v + g_{\pm}(v), \quad (4.3)$$

where  $g_{\pm}(v) = h_{\pm}(v) - \gamma v$ . The functions  $u = h_{\pm}(v)$  stand for the two branches of  $f(u) = v$ , as shown in Fig. 1. Thus, the whole plane  $R$  is decomposed into two subplanes  $\Omega_{\pm}(t)$  where  $u = h_{\pm}(v)$  hold. The boundary between  $\Omega_+(t)$  and  $\Omega_-(t)$  is an interface, say  $\Gamma(t)$ . The second stage is the motion of interfaces. As in the similar way to the scalar version, we obtain

$$\tau V = \{-c(q) + \varepsilon K\}n \quad (4.4)$$

where  $q$  is the value of  $v$  on the interface. We thus have the evolution equation for  $(v, \Gamma)$

$$\begin{cases} \frac{\partial v}{\partial t} = D\Delta v + g_{\pm}(v), & (t, x) \in \Omega_{\pm}(t) \\ \tau V = -\{c(q) + \varepsilon K\}n & \text{on } P(t) \end{cases} \quad (4.5)$$

with suitable continuity conditions for  $v$  on the interfaces. For the one dimensional case, the global existence of a smooth solution is proved in [HNM]. On the other hand, for higher dimensional cases, complicated singularities may appear so that we can not expect the global existence of smooth solutions. Recently, the local existence of a smooth solution has been shown in [Ch].

For the dynamics of interfaces, there is a great difference between the scalar case (4.1) and the system (4.4). For the former case, if the configuration of an initial-

interface is convex, it is still convex for any time, while it possibly becomes non-convex for the latter case ([OMK]).

The analysis of this system is in the beginning stage.

## 5. Concluding Remarks

In the previous sections, we have two different systems. One, system (1.5). is the so-called reaction diffusion system with a small parameter  $\varepsilon$  and the other (4.5) is the interface equation associated with (1.4) in the limit as  $\varepsilon \downarrow 0$ . The relation between (1.5) and (4.5) can be studied by intuitive and formal asymptotic analysis. The rigorous understanding has been achieved for the scalar case (3.1) ([MS]) and the numerical study ([IK]). This is one of the important problems which we should study in future, from pattern formation view point.

*Note after submission.* After completing this work, the author learned the paper by X. Chen: Generation and propagation of interfaces in reaction diffusion systems (IMA, preprint series # 708), which discuss the relation of (1.5) and (4.5).

## References

- [Ca] Caginalp, G.: An analysis of a phase field model of a free boundary. *Arch. Rat. Mech. Anal.* **92** (1986) 205–245
- [Ch] Chen, X.Y.: Dynamics of interfaces in reaction diffusion systems. *Hiroshima Math. J.* (to appear)
- [CGG] Chen, Y.G., Giga, Y., Goto, S.: Uniqueness and existence of viscosity solutions of generalized mean curvature equations. Preprint
- [EHT] Ermentrout, G.B., Hastings, S.P., Troy, W.C.: Large amplitude stationary waves in an excitable lateral-inhibitory medium. *SIAM J. Appl. Math.* **44** (1984) 1133–1149
- [ES] Evans, L.C., Spruck, J.: Motion of level sets by mean curvature I. Preprint
- [Fi1] Fife, P.C.: Understanding the patterns in the BZ reagent. *J. Statist. Phys.* **39** (1985) 687–703
- [Fi2] Fife, P.C.: Dynamics of internal layers and diffusive interfaces. *CBMS-NSF Regional Conference Series in Applied Mathematics*, vol. 53 (1988)
- [FH] Fife, P.C., Hsiao, L.: The generation and propagation of internal layers. *Nonlinear Analysis* **12** (1988) 19–41
- [FM] Fife, P., McLeod, J.B.: The approach of solutions of nonlinear diffusion equations to travelling front solutions. *Arch. Rat. Mech. Anal.* **65** (1977) 335–361
- [GM] Gierer, A., Meinhardt, H.: A theory of biological pattern formation. *Kybernetika* **12** (1972) 30–39
- [Gr] Grayson, M.A.: The heat equation shrinks embedded plane curves to round points. *J. Diff. Geom.* **26** (1987) 285–314
- [HNM] Hilhorst, D., Nishiura, Y., Mimura, M.: A free boundary problem arising from some reaction-diffusion system. Preprint
- [IMN] Ikeda, H., Mimura, M., Nishiura, Y.: Global bifurcation phenomena of traveling wave solutions for some bistable reaction-diffusion systems. *Nonlinear Analysis* **13** (1989) 507–526

- [IK] Ikeda, T., Kobayashi, R.: Numerical simulations to interfacial dynamics. 1989 (videotape)
- [KT] Keener, J.P., Tyson, J.J.: Spiral waves in the Belousov-Zhabotinsky reaction. *Physica D* **21** (1986) 307–324
- [MK] Meinhardt, H., Klinger, M.: Pattern formation by a coupled oscillations: The pigmentation patterns on the shell of molluscs. (*Lecture Notes in Biomathematics*, vol. 71.) Springer, Berlin Heidelberg New York 1987, pp. 184–198
- [MS] de Mottoni, P., Schatzman, M.: Development of interfaces in  $N$ -dimensional space. 1989 (preprint)
- [Mu] Murray, J.D.: A prepattern formation for animal coat markings. *J. Theor. Biol.* **88** (1981) 161–199
- [NAY] Nagumo, J., Arimoto, S., Yoshizawa, S.: An active pulse transmission line simulating nerve axon. *Proc. Inst. Radio Engrs.*, vol. 50, 1962, pp. 2061–2070
- [NM] Nishiura, Y., Mimura, M.: Layer oscillations in reaction-diffusion systems. *SIAM J. Appl. Math.* **49** (1989) 481–514
- [NMIF] Nishiura, Y., Mimura, M., Ikeda, H., Fujii, H.: Singular limit analysis of stability of traveling wave solutions in bistable reaction-diffusion systems. *SIAM J. Math. Anal.* **52** (1990) 142–164
- [OMK] Ohta, T., Mimura, M., Kobayashi, R.: Higher-dimensional localized patterns in excitable media. *Physica D* **34** (1989) 115–144
- [Wi] Winfree, A.T.: *The geometry of biological time*. Springer, Berlin Heidelberg New York 1980



# The Development of Algebraic Methods of Problem-Solving in Japan in the Late Seventeenth and the Early Eighteenth Centuries

Annick M. Horiuchi

REHSEIS (CNRS), 49 rue Mirabeau, F-75016 Paris, France

## I. Introduction

The rapid growth of mathematical knowledge during the Edo period (1600–1868) is one of the most remarkable features of the history of science in Japan before the modernization of the Meiji era. The chief outcomes of this long-standing tradition were:

- the accumulation of a significant body of high-standard mathematical works jealously kept within private academies where they were communicated to small numbers of selected disciples,
- the wide diffusion of mathematical practice through the publication of popular textbooks and the activity of schoolmasters and itinerant teachers.

To appreciate the extent of this development, one must recall that, at the turn of the 17th century, mathematics meant little more than elementary computations performed with the abacus. The situation was modified to a considerable extent after the Japanese turned their attention to ancient Chinese scientific works. In less than half a century, Japanese mathematics developed from the primary art of computation which served the practical needs of merchants, craftsmen and low-grade warriors into a discipline appealing to a scholarly audience.

Before dealing with the original contribution of Japanese mathematicians, it is to be noted that the assimilation of Chinese mathematics was quite effectively prompted by the novel practice of leaving several problems unsolved at the end of the books as a challenge to other mathematicians. Almost all the difficult problems discussed in Chinese works were integrated in this way into the Japanese corpus.

A further impetus was added in 1658 by the discovery and the subsequent reprint of a 13th century Chinese treatise [Zhu Shijie 1299], the level of which far surpassed those of previously available works. The 13th century in the history of Chinese mathematics was a very productive period (often described as its golden age) when significant achievements were made, most particularly in the area of algebraic devices [Li Yan and Du Shiran 1987, chap. 5]. The Japanese scholars' attention soon focused on the *tengen* (or *tianyuan* in Chinese) method which Zhu Shijie used to solve a range of tricky problems.

The *tengen/tianyuan* method of solving problems was basically similar to the one which is called algebra in the West. The Japanese scholars spent many years before getting at the meaning of the word *tengen* (literally “celestial origin”), the name given by the Chinese to the “unknown”. The level of the difficulties involved was such that the contemporary Chinese mathematicians who were unaware of the past achievements had to wait for the introduction of the western algebra in the 18th century to rediscover the meaning of the method. [Martzloff 1987, pp. 105–106]

Japanese mathematicians worked on problems which were modelled on the older Chinese problems. These problems dealt with concrete situations and were expressed in numerical terms; their solution included both the numerical result and the procedure (*jutsu* in Japanese), the sequence of the operations to perform on the abacus or with counting-rods to get to the result.

The proliferation of small problems of this kind throughout the Edo period and the increasingly artificial character of most of them led some historians to stress the artistic and recreational character of Japanese mathematics and the mathematicians’ indifference to the ‘utility’ of their art [Mikami, 1921].

This point of view relied chiefly on the examination of one side of Japanese mathematics: the problems. But one cannot ignore the fact that the “utility” of mathematics in the past very often depended on the general and efficient *methods* and *tools* which were obtained through solving particular problems. Therefore, the issue of the utility of Japanese mathematics cannot be settled before having examined the other side of the mathematicians’ work, that of elaborating *methods* of solving problems.

The aim of this paper is to discuss some prominent features of the development of these methods in the late 17th and the early 18th centuries. I will focus on the achievements of Seki Takakazu (?–1708) and Takebe Katahiro (1664–1739), two major mathematicians of this period. I will discuss the way Seki improved the Chinese algebraic methods and stress the importance of the Chinese root-extraction procedure in the course of his research. I will then turn to one of Takebe’s main contributions, the introduction of the infinite power series in the scope of algebraic calculation.

## II. SEKI’S STUDY OF ALGEBRAIC DEVICES

Let us begin with Seki’s synthesis of earlier methods of problem-solving. The synthesis, achieved by 1683, took the form of a trilogy. Each treatise was devoted to a particular method of problem-solving [Hirayama et al. 1974, sects. 6, 7 and 8]. I will examine only the last treatise where Seki expounds an original method of problem-solving which can be understood as an extension of the Chinese *tengen/tianyuan* method.

Let us consider first the main features of this Chinese method which had a considerable influence on the development of algebraic devices in Japan. Here is an example of a problem requiring the *tengen* method for its solution:

“Given a rectangular field of 8 mu 5 fen and 5 li [1 mu = 240 square bu; 1 fen = 0.1 mu; 1 li = 0.1 fen]. We only say that the sum of the length and the width is 92 bu. Find the length and the width”. [Zhu Shijie 1299, chapter kaifang shisuo]

The solution followed a regular pattern: first, a quantity was set up as the unknown. Then two different algebraic expressions for a certain quantity were built up (in the problem above, the area was expressed as  $x(92 - x)$  and 2052). The equation was derived by subtracting one of the expressions from the other. Finally, the numerical value of the unknown was determined by extracting the root of the equation, digit by digit.

One basic feature of the *tengen/tianyuan* method was the use of counting rods to carry out the polynomial calculation as well as the extraction of the root. The solution itself, as I have said in the introduction, consisted of the sequence of operations to be performed with this instrument. The instructions given by the author were like: "Put one rod for the unknown", "Multiply by itself four times", "Add it to the area", etc.

The way the counting-rods were to be handled could be reconstructed from specific symbols representing the polynomials obtained on the counting-board (a large sheet of paper with horizontal and vertical lines drawn on it). These symbols were inserted in the solution after each instruction. Polynomial expressions as well as equations of one unknown were represented by the column of their coefficients arranged in the order of increasing powers of  $x$  (see Fig. 1):

	209 (constant term)
	-16 (coefficient of $x$ )
	3 (coefficient of $x^2$ )
Ex: $3x^2 - 16x + 209$	

Fig. 1. The Chinese notation for polynomials

Let us now turn to Seki's improvement of the Chinese method. The need for an improvement originated in the publication in 1671 of fifteen unsolved problems by Sawaguchi Kazuyuki, a Kyoto mathematician. Sawaguchi's problems were so intricate that none of them could be handled with the Chinese method. This is an example of Sawaguchi's problems:

"We have now  $A$ ,  $B$  and  $C$ , such that each is a cube.

We say first that the volumes of  $A$  and  $B$  altogether make 137,340 tsubo and also that the volumes of  $B$  and  $C$  altogether make 121,750 tsubo.

We say in addition that the square root of the edge of  $A$ , the cube root of the edge of  $B$  and the fourth root of the edge of  $C$  altogether make 1 shaku 2 sun. Find the sides of  $A$ ,  $B$  and  $C$ ." [Sawaguchi 1671, Problem 4].

The third treatise of Seki's trilogy gave a global solution to two questions implicitly raised by Sawaguchi's problems which can be formulated as follows: How can one proceed to the calculation when the quantities involved cannot be

	$(315l - 4h)$ (constant term)
	$(-12h^5)$ (coefficient of $x$ )
	$3bhl$ (coefficient of $x^2$ )

$$\text{Ex: } 3bhlx^2 - 12h^5x + (315l - 4h)$$

Fig. 2. Seki's notation for polynomials with literal coefficients.

The literal and the numerical parts were dealt with separately. The latter part was transcribed by using the rod-numerals; the former was written beside these numerals using Chinese characters.

In the example above, the Chinese characters meaning length, breadth, height and number four have been replaced by the letters  $l$ ,  $b$  and  $h$  and the arabic numeral 4.  $h^5$  is represented by  $\frac{h}{4}$  because Japanese mathematicians used to consider the number of times a quantity was multiplied by itself, that is the power minus one unit.

expressed in terms of one unknown? How can one eliminate the unknown within two equations?<sup>1</sup>

Seki answered by adding supplementary steps to the *tengen* pattern. In the course of these steps, additional unknowns were introduced and eliminated. Seki's improvement can be characterized by two main features. First, the general pattern of the *tengen* was maintained by considering one unknown at each stage and by integrating the other unknown quantities into the data. Second, the calculation was reduced to a single process of eliminating one unknown within two given algebraic equations.

Seki's ability to cope with such general contexts and questions was closely related to his use of adequate notations to represent polynomials and equations with literal coefficients (see Fig. 2). Seki's notations were obtained by extending the traditional representation of polynomials with numerical coefficients. The rules of calculation with the new notations remained unchanged. [Hirayama et al. 1974, Sect. 29].

This step has been described by many historians as a shift from an "instrumental" algebra into a "written" algebra. But this description is misleading in at least one respect. In fact, the solution maintained its algorithmic character and consisted of a rhetorical description of the operational instructions to be performed on an imaginary counting-board. In the solution, the notation was only used to represent the polynomial expressions obtained at each main step of the calculation. The calculation itself was performed somewhere else and was not made explicit. We do not even know if Seki performed the calculation on the paper or went on using the counting-board in some manner. In this respect, his calculation still preserved an instrumental character.

<sup>1</sup> Regarding this point, an extension to the case of four unknowns was achieved in China as early as the beginning of the 14th century [Zhu Shijie, 1303]. Japanese mathematicians of Edo period did not know about this extension.

Additional features should be noted about Seki's use of this notation. Seki described the procedure of elimination in very general terms, by considering sets of two equations with arbitrary coefficients.

The coefficients were represented by means of Chinese characters taken from a series of twelve characters (the ten “stems” and the twelve “branches”; *kanshi*) in the same way as alphabet letters symbolise numbers in western algebra.

Many studies have been devoted to Seki's method of elimination by which he solved problems through cancelling a certain determinant [Mikami 1913]. Leaving aside detailed analysis of this method, we note that Seki's interest in the general rules of calculating the determinant was closely connected to his systematic use of Chinese characters as symbols in place of particular numbers.

To conclude the first part of this study, I will add that the whole trilogy of Seki's reveals a clear shift of his concern from particular problems to general methods of solution. In all the treatises composing the trilogy, Seki displays both his intention and his ability to reduce any particular problem to general processes of calculation.

This tendency towards generalization, however, was not entirely new in the Sino-Japanese tradition. The search for general procedures of problem-solving had always been part of the Chinese tradition, as is well exemplified by the classification adopted in one of the oldest mathematical treatises in China, the *Nine Chapters of Mathematics* [Chemla 1988; Wu 1986]. Seki's originality lay rather in departing from the general tendency of his time to favour particular problems, and in his introduction of a range of efficient tools to describe the methods in general terms. As a consequence, Seki's attention gradually concentrated on general objects like equations and polynomials.

### III. ‘Defective’ Problems and Seki’s Reflection on Root-Extraction Procedures

Let us now examine more closely Seki's study of equations. This part of Seki's achievement, which took place shortly after the trilogy, is particularly interesting in that it brings forward a quite different way of discussing properties of equations and roots from the one so far better-known to modern mathematicians. A particular computational device plays a central role in Seki's research: the root-extraction procedure.

Seki's interest in equations stemmed from a particular, quite pragmatic preoccupation: the search for methods to correct what he called ‘defective’ problems (*byôdai*). According to Seki's definition [Hirayama et al. 1976, Sect. 9], problems were “defective” or wrongly stated if they led either to more than one acceptable solution or to none. Defective problems had to be corrected by changing the terms of the problem.

In the course of his research, Seki's interest shifted to the equations themselves and he thought out a general method of transforming an equation with more than one root into one equation with a single root [Hirayama et al. 1974, Sect. 8]. This was done by changing the value of one of the coefficients in the initial equation.

Before discussing Seki's method, we must note that the notion of "equation" did not have a strict equivalent in traditional Chinese mathematics. Instead, we find the concept of "configuration for extraction" (*kaifangshi*), referring to the numbers set up on the counting-board to perform the extraction. The configuration itself was similar to the polynomial expression (Fig. 1). Likewise, instead of the concept of "root", we find that of "quotient", which referred to the location on the counting-board of the result of the extraction. The coefficients of the equation were similarly called by the names of their respective locations on the counting-board.

The extraction procedure which Seki used in this context was a general procedure which allowed him to compute successively all the real roots of a given equation.

This procedure, in fact, was an outcome of a long tradition of research in China [Li Yan and Du Shiran 1987, Sect. 5.2] but let us concentrate here only on the procedure as was set forth in Seki's treatise [Hirayama et al. 1976, Sect. 8]. The whole extraction was built on one basic pattern of computation involving the coefficients of the equation (see Table 1), in which one can easily recognize the so-called Horner-Ruffini process. The successive configurations can be interpreted as a gradual alteration of the coefficients of (1) into the coefficients of (2) where (2) is satisfied by  $y$  with  $y = x - a$ .

$$f(x) = m + nx + lx^2 + px^3 = 0 \quad (1)$$

$$\varphi(y) = f(y + a) = 0 \quad (2)$$

The above interpretation of the root-extraction pattern as a substitution of the unknown was not explicitly stated by Seki in this context. But all the improvements he introduced suggest that Seki did have a similar explanation of it.

Let us examine more closely how Seki used this procedure to find all the roots of a given equation. To begin, the first root of the equation, let us call it  $a$ , was sought by carrying out successively this basic pattern with several quotients

**Table 1.** The basic pattern of Seki's root extraction procedure

(1)		(2)	
Quotient	$a$	$a$	
Constant term	$m$	$m + (n + la + pa^2) \cdot a = m + na + la^2 + pa^3$	
Coefficient of $x$	$n$	$n + (l + pa) \cdot a = n + la + pa^2$	
Coefficient of $x^2$	$l$	$l + p \cdot a$	
Coefficient of $x^3$	$p$	$p$	
(3)		(4)	
$a$		$a$	
$m + na + la^2 + pa^3$		$m + na + la^2 + pa^3$	
$n + la + pa^2 + (l + 2pa) \cdot a = n + 2la + 3pa^2$		$n + 2la + 3pa^2$	
$l + pa + p \cdot a = l + 2pa$		$l + 2pa + p \cdot a = l + 3pa$	
$p$		$p$	

The calculation is carried out for each configuration from the bottom up.

(positive or negative), suitably chosen so that the number in the top line gets closer and finally becomes equal to zero. The first root was then derived by adding the list of quotients.

The equation  $\varphi(y) = 0$  corresponding to this last configuration would therefore be inferior by one unit compared to the initial equation.

Seki's originality lay in his idea of iterating the previous process to the new equation. Assuming for example that  $b$  is a root of the new equation  $\varphi(y) = 0$ , which means that at the end of the iterated process with  $b$  as quotient, the top line of the last configuration gets to zero, then the number  $a + b$  would be the second root of the initial equation.

The central role played by this root-extraction procedure in Seki's reflection on equations and roots is particularly obvious in the way he tried to solve the aforementioned question of removing the "excess" roots.

Seki's method was based on a general criterion which was to be fulfilled by the coefficients of any equation having a single root. This criterion which is actually a condition of existence of a double root was the outcome of a computation in which the extraction procedure played a central role.

The fact that  $a$  is the single root of (1) meant for Seki that at the end of the extraction procedure, the top two lines of the last configuration (the configuration (4) in Table 1) were equal to zero.

$$\begin{cases} m + na + la^2 + pa^3 = 0 \\ n + 2la + 3pa^2 = 0. \end{cases}$$

The criterion was then derived by eliminating  $a$  within these two equations:

$$27m^2p^2 + 4ml^3 + 4n^3p - 18mnlp - n^2l^2 = 0.$$

As shown by this example, the root-extraction procedure in Seki's mathematics was clearly something more than a device for computing the roots of an equation. This procedure was automatically involved in any discussion pertaining to equations, coefficients or roots. No hypothesis on the roots or any property of the equation could be stated without referring to this procedure. "Theory of extraction" would thus be the most suitable name for this part of Seki's research.

#### IV. Takebe's Method of Computing the Length of An Arc

Let us turn now to another significant response given by the Japanese mathematicians to a very ancient problem of trigonometry, the problem of finding a general procedure for calculating the length of an arc  $a$  in terms of the sagitta  $s$  and the diameter of the circle  $d$  (see Fig. 3).

We shall focus here on the work of Seki's disciple, Takebe Katahiro, whose contribution provides the crowning step of a long tradition of research in China and Japan [Li Yan and Du Shuran, § 3.2]. In 1720, Takebe expressed for the first time the square of the length of the arc of the circle as an infinite power series of the sagitta. Takebe's work displays an unprecedented intensity in the application of algebraic methods and especially of the root-extraction procedure.

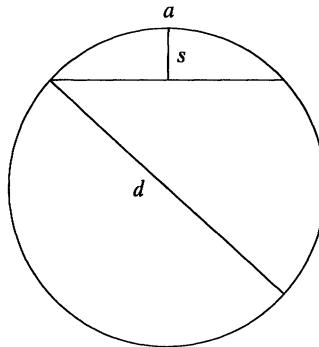


Fig. 3

To begin with, let us look at the infinite procedure of calculating the quantity  $(a/2)^2$  as formulated in Takebe's treatise [Takebe 1722, sect. 12].  $a$ ,  $s$  and  $d$  are respectively the length of the arc, the sagitta and the diameter of the circle.

“The fundamental procedure runs as follows:

Sagitta and diameter are multiplied. This gives the approximate value of the square of half the arc. [ $sd$ ]

Divide by three the square of the sagitta. This gives the first discrepancy.  
[ $X_1 = s^2/3$ ]

Put down the first discrepancy and multiply by the sagitta. Divide by the diameter. Then, multiply by 8 and divide by 15. This gives the second discrepancy. [ $X_2 = X_1 \cdot (s/d) \cdot (8/15)$ ]

Put down the second discrepancy. Multiply by the sagitta. Divide by the diameter. Multiply by 9. Divide by 14. This gives the third discrepancy. [ $X_3 = X_2 \cdot (s/d) \cdot (9/14)$ ]

Put down the third discrepancy. Multiply by the sagitta. Divide by the diameter. Multiply by 32. Divide by 45. This gives the fourth discrepancy.  
[ $X_4 = X_3 \cdot (s/d) \cdot (32/45)$ ]

[...]

The successive discrepancies are added to the approximate value of the square of half the arc. This gives the fixed value of the square of half the arc.”  
[Takebe 1722, Chap. 12]

Leaving aside the highly empirical method through which Takebe found his way to this procedure, we shall focus on his extension of the scope of algebraic calculation.

The first thing to examine is the original commentary he added to the explanation of his method. In this commentary, Takebe took as his starting point a common distinction between ‘exhaustible’ numbers (*tsukuru kazu*) and ‘inexhaustible’ ones (*tsukizaru kazu*). Numbers were exhaustible (respectively inexhaustible) if they had a limited (respectively unlimited) decimal expression. He then extended this distinction to procedures and formed the following programme:

"Numbers like  $1/4$  and  $1/5$  are exhaustible numbers.  $1/3$  and  $1/7$  are examples of inexhaustible numbers. Addition, subtraction, and multiplication are exhaustible procedures. Division and extraction are inexhaustible procedures. The perimeter of the square or the areas of rectangles have an exhaustible form. The circumference and the segment of the circle have an inexhaustible form. Thus, just as the forms of the arc and the circle are inexhaustible, the procedure involved is also inexhaustible. Since the procedure is inexhaustible, so are the resulting numbers."

[Takebe 1722, Chap. 12]

Takebe's commentary aimed at conferring a legitimate status to the infinite procedure, a thing which must have been quite unfamiliar and uncomfortable to most of his contemporaries. The link he established between the infinite procedures and the numbers with an infinite decimal expression may be considered as his primary argument for integrating the new procedure into the existing mathematical corpus.

The role played by this link was not only pedagogical. As can be seen in Takebe's later works, this analogy between the decimal numbers and the infinite series allowed him to venture into still more new fields of research.

This last point is exemplified by the second method of calculating the length of arc which Takebe devised shortly after the first. The method rested on the idea of deriving an infinite series by performing the root-extraction on an equation with literal coefficients.

We shall limit our discussion to the first step of the process which contains the main idea. Takebe resorted to a classical device consisting in inscribing polygonal lines inside the arc and getting approximate values of the arc by calculating the length of the polygonal lines (see Fig. 4).

Now, let us call  $a_0$  the length of the arc to be calculated and  $s_0$ , the sagitta in terms of which the arc is to be expressed.

As can easily be seen, the sagitta  $s_1$  of the arc  $a_1 = a_0/2$  satisfies the following equation (1).

$$-s_0 d + 4dx - 4x^2 = 0 \quad (1)$$

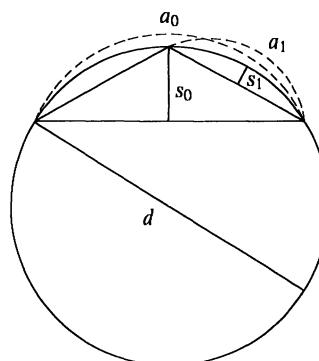


Fig. 4

The extraction procedure is performed on (1) exactly as if the coefficients were numerical. As a result of this extraction  $s_1$  is expressed as an infinite power series of  $s_0/d$ .

$$s_1 = s_0/4 + s_0^2/16d + s_0^3/32d^2 + 5s_0^4/256d^3 + 7s_0^5/512d^4 + 21s_0^6/2048d^5 + \dots$$

The same process is then performed on (2) in which the constant term  $-s_1d$  is given as an infinite series of  $s_0/d$ .

$$-s_1d + 4dx - 4x^2 = 0 \quad (2)$$

$s_2$  could then be expressed in terms of  $s_0/d$  by performing the root-extraction procedure on (2). We note that in the course of the calculation successive terms of the infinite series are handled as if they were successive digits of a decimal number.

This example shows clearly the central role played by the aforescribed analogy in Takebe's second method of calculating the arc. His extension of the extraction procedure to cases where coefficients are literal and even infinite series, as well as his way of handling the infinite series, were direct outcomes of this analogy.

From a historical point of view, this later aspect of Takebe's research had a fundamental significance in offering a very effective method of constructing infinite series. Moreover, in his explanation of the second method, Takebe clearly displayed his belief that the previously described extraction procedure could be applied to an equation of any degree. He also suggested that the general law which governs the successive terms of the series was resulting from the extraction procedure [Takebe s.d, p. 24b]. Both points bring Takebe's method very close to the idea of the binomial theorem.

## Conclusion

Even though the analyses above do not cover the whole scope of mathematical research in the Edo period, they show well enough the growing interest in algebraic devices during the 17th and the 18th centuries in Japan. We have seen how the Chinese algebraic devices, including the procedure of extraction, have been extended and enriched by Seki and Takebe. By studying these devices, they not only sought to solve the largest amount of problems, but also strove to clarify and make explicit the general patterns of methods of problem-solving. In my view, this outstanding feature of the Japanese tradition has more claim to the utilitarian cause than has usually been recognized.

## References

- Chemla, Karine (1988): La pertinence du concept de classification pour l'analyse des textes mathématiques chinois. Extrême-orient, Extrême-occident, 10
- Hirayama, Akira, Shimodaira, Kazuo, Hirose, Hideo (eds.) (1974): Seki Takakazu zenshû (Seki Takakazu's Complete Works). Osaka kyôiku tosho, Tokyo
- Lam, Lay Yong (1982): Chinese polynomial equations in the thirteenth century. In: Explorations in the history of science and technology in China. Shanghai

- Li, Yan, Du, Shiran (1987): Chinese mathematics – A concise history. Clarendon Press, Oxford
- Martzloff, Jean-Claude (1987): Histoire des mathématiques chinoises. Masson, Paris
- Mikami, Yoshio (1913): On the Japanese theory of determinants. *Isis* 1, 9–36
- Mikami, Yoshio (1921): *Bunkashijō yori mitaru nihon no sūgaku* (Japanese mathematics considered from the point of view of cultural history). Kōseisha Kōseikaku, Tokyo (reprint, 1984)
- Qin, Jiushao (1247): *Shushu jiuzhang* (Mathematical treatise in nine chapters)
- Sawaguchi, Kazuyuki (1671): *Kokon sanpōki* (Treatise of ancient and modern mathematics)
- Takebe, Katahiro (1722): *Tetsujutsu sankyō* (The mathematical classic of *tetsujutsu*)
- Takebe, Katahiro: *Enri kohai jutsu* (The procedure for calculating the arc based on the principle of the circle)
- Wu, Wentsun (1986): Recent studies of the history of Chinese mathematics. Proceedings of the International Congress of Mathematicians, Berkeley, California, USA. Amer. Math. Soc., Berkeley 1987
- Zhu, Shijie (1299): *Suanxue qimeng* (Introduction to mathematical studies)
- Zhu, Shijie (1303): *Siyuan Yujian* (The jade mirror of four unknowns)



# The Birth of Spectral Theory – Joseph Liouville’s Contributions

*Jesper Lützen*

Mathematics Department, University of Copenhagen, Universitetsparken 5  
DK-2100 Copenhagen Ø, Denmark

## 1. Introduction

The history of spectral theory is the history of a beautiful and important area of mathematics with close links to physics and with a strong influence on the development of functional analysis. Its roots lies in three areas: 1) discrete systems described by matrices (or quadratic forms) and continuous systems described by 2) differential equations or 3) integral equations. These different appearances of spectral theory were not formally connected until around 1900. The history of spectral theory of matrices has been studied in detail by Hawkins (e.g. 1975), so here I shall concentrate on the last two areas. In both of these, Joseph Liouville (1809–1882) played a role in the early period. His and his friend Charles Sturm’s (1803–1855) work from the 1830s on spectral theory of self adjoint second order differential equations is emphasized in all histories of mathematics of the 19th century (e.g. Dieudonné 1981). The bulk of his work on spectral theory of integral operators, on the other hand, remained unpublished in his notebooks, until I reconstructed and published it in a scientific biography of Liouville that has recently appeared (Lützen 1990 Chap. XV). Therefore Fredholm’s and Hilbert’s works from the beginning of the 19th century are the earliest published works on spectral theory of integral operators, but Liouville’s notes from the mid 1840s show that he anticipated many of the basic ideas with more than half a century. Moreover they reveal that Liouville used a variational technique, named after Rayleigh and Ritz (1877 and 1909 respectively), to determine the eigenfunctions<sup>1</sup>, and that he questioned the naive version of the Dirichlet principle more than 20 years before Weierstrass’ famous criticism.

Liouville’s “integral operator” originated in a study of potential theory of charged surfaces, and Liouville’s approach has turned out to be of interest even to modern potential theorists (Berg and Lützen 1990), (Berg and Fuglede 1990), (Bang and Fuglede 1990), (Berg 1990).

---

<sup>1</sup> In the 1810s continuous spectra had implicitly been considered by Fourier and Cauchy in their work on the Fourier Integral. Otherwise continuous spectra do not occur until 1897 when Wirtinger was led to them in his study of Hill’s equation (Dieudonné 1981 p. 149). Therefore, in this paper, we shall only consider discrete spectra of eigenvalues.

## 2. Separation of Variables

Sturm and Liouville considered the following differential equation

$$(k(x)V'(x))' + (g(x)r - l(x))V(x) = 0 \quad \text{for } x \in (\alpha, \beta) \quad (1)$$

with the boundary conditions

$$k(x)V'(x) - hV(x) = 0 \quad \text{for } x = \alpha \quad (2)$$

and

$$k(x)V'(x) + HV(x) = 0 \quad \text{for } x = \beta \quad (3)$$

where  $k, g$  and  $l$  are positive functions,  $h$  and  $H$  are positive constants and  $r$  is a parameter. They both arrived at this problem by separating variables (i.e. by setting  $u(x, t) = V(x)e^{-rt}$ ) in the partial differential equation

$$g \frac{\partial u}{\partial t} = \frac{\partial(k \frac{\partial u}{\partial x})}{\partial x} - lu \quad (4)$$

governing the heat conduction in a heterogeneous, unequally polished bar.<sup>2</sup> Here  $u$  represents the temperature, and the boundary conditions (2) and (3) reflect that the ends of the bar are maintained at zero degrees.

Sturm and Liouville saw that there were only non-trivial solutions to (1)–(3) when  $r$  belongs to a countable set of eigenvalues<sup>3</sup>  $r_1, r_2, r_3, \dots$ , and they studied the following questions:

- 1) properties of the eigenvalues;
- 2) behaviour of the corresponding eigenfunctions  $V_1, V_2, V_3, \dots$ ;
- 3) expansion of arbitrary functions in a series of eigenfunctions.

The expansion problem arose in attempting to fit the solution  $\sum A_n V_n e^{-r_n t}$ , found by superposing simple solutions of (4), to a given initial temperature distribution  $u(x, 0) = f(x)$ . Indeed this poses the problem of finding  $A_n$ 's so that

$$\sum A_n V_n(x) = f(x). \quad (5)$$

Separation of variables was at the origin of all the early eigenvalue problems of differential operators (Lützen 1987). It was used in its full generality by Fourier in his influential work on heat conduction in homogeneous equally polished materials (1822). Stationary heat conduction in a rectangular plate and non-stationary heat conduction in a rod led Fourier to the following simple special case of (1):

$$V''(x) = -m^2 V(x) \quad (6)$$

with suitable boundary conditions.

<sup>2</sup> In his earliest work, Liouville (1830) had  $g$  and  $k$  constant.

<sup>3</sup> I shall use Hilbert's terminology "eigenvalue" "spectrum" etc. Moreover, I shall use the word "orthogonal", although Liouville had no such geometric ideas.

In this case the series (5) is the usual trigonometric Fourier series. When studying stationary heat conduction in spheres and cylinders Fourier used spherical and cylindrical coordinates in order to have one of the coordinates constant at the boundary. In these coordinates, the Laplacian has variable coefficients and so, separation of variables led to more complicated differential operators with variable coefficients, for which the eigenfunctions could only be found as infinite series.

Similar ideas had been used already in 1759 by Euler, to study vibrations of a circular membrane (eigenfunctions are Bessel functions) and implicitly in the 1780s by Legendre and Laplace in their study of attraction from a spheroid (eigenfunctions are spherical harmonics). In the 1830s Lamé developed such applications of curvilinear coordinate systems into a flexible method. In particular he studied stationary heat conduction in an ellipsoid and was thus led to the so-called Lamé functions or ellipsoidal harmonics. Later in the century the method of separation of variables applied to the partial differential equations of physically interesting phenomena gave rise to many other special functions.

### 3. Sturm-Liouville Theory

In their study of the eigenvalue problem (1)–(3) with variable, and not even explicitly given coefficients, Sturm and Liouville hit on problems of an entirely new nature. Indeed, since no workable explicit expression of the solutions can be found, they had to work directly with the equation, and the results they found were necessarily qualitative. Except for Cauchy’s theorem on the existence of a solution to the Cauchy problem of a first order differential equation, (1824–1881) all previous works on differential equations had asked the question: given a differential equation, find its solutions. Sturm and Liouville broadened the question to: given a differential equation, find some property of its solutions. Conceptually this is a great step toward a qualitative theory, taken up by Poincaré in the 1880s in the case of non-linear equations.

Sturm mainly studied point 1) and 2) above. By a variational technique, he followed the behaviour of the roots of a solution  $V_r$  of (1) and (2). He showed that these roots are decreasing functions of  $r$ . From this observation he concluded that there is a countable infinity of eigenvalues of the system (1)–(3) and he deduced his famous comparison and oscillation theorems. He developed most of his remarkable theory in the period from 1829 to 1833, but it did not appear until Liouville invited him to publish it in the first volume of his journal (Sturm 1836a,b).

#### 3.1 Liouville’s 1830 Paper

Liouville primarily studied the expansion problem 3), finding, in the process, additional results related to 1) and 2). His first paper of 1830 was a peculiar mixture of ingenious ideas and unrigorous methods. Among the ingenious ideas was his use of the method of successive approximations to express the solution of (1) and (2). He proved the convergence of the series and had thus published

the first theorem of existence of a differential equation (1830, 1836).<sup>4</sup> Moreover he established the orthogonality relation

$$\int_{\alpha}^{\beta} g(x) V_m(x) V_n(x) dx = 0 \quad \text{for } m \neq n. \quad (7)$$

This is the only general result in Sturm-Liouville theory, that had been shown earlier, namely by Poisson (1826).

Multiplying (5) by  $g(x)V_m(x)$  and integrating from  $\alpha$  to  $\beta$ , Liouville found by (7), that if  $f(x)$  can be expressed in a series eigenfunctions, the “Fourier” coefficients must be of the form

$$A_n = \frac{\int_{\alpha}^{\beta} g(x) V_n(x) f(x) dx}{\int_{\alpha}^{\beta} g(x) V_n^2(x) dx}. \quad (8)$$

In this argument Liouville integrated term by term in an infinite series. We know that this is problematic, but he never questioned this exchange of limit procedures.

Finally Liouville determined the asymptotic behaviour of  $V_n$ :

$$V_n(x) \sim \frac{\sin \sqrt{n} x}{\sqrt{n}} \quad \text{for } n \text{ large,} \quad (9)$$

and used this to prove the convergence of the “Fourier” series of  $f$ :

$$\sum_n A_n V_n = \sum_n \frac{\int_{\alpha}^{\beta} g(x) V_n(x) f(x) dx}{\int_{\alpha}^{\beta} g(x) V_n^2(x) dx} V_n. \quad (10)$$

In 1830 he used the expression for  $V_n$ , found by successive approximations, to establish the latter two theorems, and his proofs are highly questionable even with the standards of rigour of that time.

### 3.2 Liouville's Mature Papers

Liouville returned to these questions in a series of papers from the period 1836–1837. Now he did much better, for two reasons; 1) he could use Sturm's investigations of the oscillatory behaviour of  $V_n$  and 2) he could refer to Dirichlet's proof (1829) of the convergence of ordinary trigonometric Fourier series. With these tools Liouville gave a beautiful and rigorous proof (in fact two of increasing generality) (1837a,b) of the convergence of the Fourier series (10) for a large class of functions  $f$ . This may be considered his most important contribution to Sturm-Liouville theory.

However the problem remains: what is the limit of the Fourier series? In his first paper in the series of 1836–37 (1836), Liouville claimed that the Fourier series converges to the function  $f$ , i.e. if

---

<sup>4</sup> Cauchy's first existence proof was not published until 1981. The method of successive approximations was later attributed to Picard (1890).

$$F(x) = \sum_{n=1}^{\infty} \frac{\int_{\alpha}^{\beta} g(x) V_n(x) f(x) dx}{\int_{\alpha}^{\beta} g(x) V_n^2(x) dx} V_n(x) \quad (11)$$

then  $F = f$ . In fact he showed that

$$\int_{\alpha}^{\beta} g(x)(F(x) - f(x)) V_n(x) dx = 0 \quad \text{for all } n \quad (12)$$

and inferred that  $F(x) - f(x) = 0$ . As Liouville himself later realized, the last inference is wrong. One can only conclude that  $F - f$  oscillates “infinitely fast”. In a note in his notebooks he attributed this insight to a certain Mr. D, probably his good friend Dirichlet, who had explicitly excluded such functions from his convergence proof. Despite repeated efforts by Liouville, partly in collaboration with Sturm, this major lacuna in Sturm-Liouville theory was left open until around 1900.

### 3.3 Higher Order Sturm-Liouville Theory

In a course at the College de France, Liouville (1838) generalized Sturm-Liouville theory to certain non-self adjoint higher order equations. Already Lagrange had introduced the adjoint equation and now Liouville introduced the adjoint boundary values. With this new concept he could prove a biorthogonality relation, that takes the place of (7) and he succeeded in generalizing most of Sturm’s results. However, he was not able to generalize his own main theorem on the convergence of Fourier series, in spite of the importance he attached to this question. There were good reasons for this failure: Even in simple cases the Fourier series does not converge to the expanded function (see Haagerup’s example in Lützen 1984).

Why did Liouville’s intuition fail in this case? I think the reason is that Liouville (as well as many of his French contemporaries) were often inspired by physics and borrowed their intuition from this science. However the generalizations proposed by Liouville did not correspond to any physical problem. Thus ended the first epoch in the history of Sturm-Liouville theory with an only partially successful generalization for the sake of generality. When the theory was revived about 40 years later, it was again a physically interesting problem that lay behind, and this was also the case when Liouville in the 1840s took up spectral theory of integral operators.

## 4 Liouville on Integral Operators

Fourier and Laplace transforms had given rise to integral equations already in the 18th century, but if we disregard these implicit occurrences, Abel (1823) was the first to study such an equation. Liouville broadened the approach in the 1830s by emphasizing the great physical interest of integral equations (1832) and by offering his theory of differentiation of arbitrary order, as a means to solve them. Moreover integral equations played a central role in his work on Sturm-Liouville theory (Lützen 1982).

#### 4.1 The Published Note

At the end of 1845 Liouville published a  $1\frac{1}{2}$  page note on spectral theory for a general class of integral operators. He considered the eigenvalue equation<sup>5</sup>

$$\int_D l(x') T(x, x') \zeta(x') dx' = m \zeta(x) \quad (13)$$

where  $l$  is a real valued function defined in a subset  $D$  of  $\mathbb{R}^n$  (or perhaps  $D$  is an  $n$  dimensional manifold) and  $T$  is real and symmetric on  $D \times D$ . His note contained two theorems:

**Theorem 1.** *Let  $\zeta_1, \zeta_2$  be solutions (eigenfunctions) corresponding to two different (eigen) values  $m_1$  and  $m_2$ . Then the following (orthogonality) relation holds:*

$$\int_D l(x) \zeta_1(x) \zeta_2(x) dx = 0 \quad \text{for } m_1 \neq m_2. \quad (14)$$

**Theorem 2.** *If  $l$  is always positive all the eigenvalues  $m$  are real.*

Liouville even concluded his brief note with the following prophetic remark:

Moreover, one can easily see that instead of the left-hand side of the equation  $A[(13)]$  one can substitute more complicated integrations or even operations of another kind without the theorem (conveniently modified, if necessary, to suit these new operations) and even the proof ceasing to be correct; for the statements above rely essentially on a certain symmetry which it will be sufficient to retain (Liouville 1845b).

With our present day knowledge of symmetric operators we must admire Liouville's foresight.

#### 4.2 The Unpublished Notes

My admiration grew considerably when, two years ago, I discovered that this note contained but a small fraction of a great spectral theory of an integral operator in potential theory that Liouville confided to his notebooks around 1846.

Liouville's notes begin as follows:

- 1°. Render  $\iiint \frac{\lambda \lambda' d\omega d\omega'}{A}$  a minimum with  $\iint \lambda d\omega = \text{const.}$ ; and you have a function  $l$  for which  $\iint \frac{l' d\omega'}{A} = \text{const.} = 1$ , say; that corresponds to the equilibrium distribution of electricity.
- 2°. Render maximum

Let  $\iint \lambda d\omega = 0$  and  $\iint \frac{\lambda^2 d\omega}{l} = \text{const.}$  and you find a second function  $l\zeta_1$  such that

$$\iint \frac{l'\zeta'_1 d\omega'}{A} = m_1 \zeta_1; \quad m_1 < 1 \quad \text{if} \quad \iint \frac{l' d\omega'}{A} = 1.$$

<sup>5</sup> Of course Liouville did not use vector notation. He wrote  $x, y, \dots, z$  instead of  $x$ .

3°. Let  $\iint \lambda d\omega = 0$ ,  $\iint \lambda \zeta_1 d\omega = 0$  and  $\iint \frac{\lambda^2 d\omega}{\Delta} = \text{const}$ . You find the function  $l\zeta_2$  such that

$$\iint \frac{l\zeta_2' d\omega'}{\Delta} = m_2 \zeta_2, \quad m_2 < m_1$$

and so on.

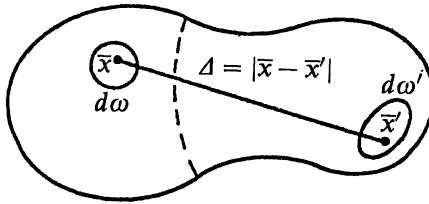
And a function  $Q$  can be expressed as

$$Q = l(A_0 + A_1 \zeta_1 + A_2 \zeta_2 + \cdots + A_n \zeta_n + \cdots),$$

$$\iint l \zeta_m \zeta_n d\omega = 0, \quad A_n = \frac{\iint Q \zeta_n d\omega}{\iint l \zeta_n^2 d\omega}.$$

That is what we found a long time ago (Liouville Ms 3618 (1), pp. 15v–16r).

What goes on here? From the rest of Liouville's approximately 100 pages of messy notes on this subject it appears that he considers a closed surface  $S$  with a single layer charge distribution of which the charge density at a point  $x$  is given by the function  $\lambda = \lambda(x)$ . If this function is considered a function of another point  $x'$  on  $S$  Liouville writes  $\lambda' = \lambda(x')$ . Moreover  $\Delta$  denotes the distance  $|x - x'|$  between the two points and  $d\omega$  is the surface measure on  $S$ .



In the first step Liouville minimizes the potential energy of the charge distribution under the assumption that the total charge is kept constant. This is similar to the way Gauss (1840) had shown the existence of the equilibrium distribution, and Liouville concludes that the minimizing function  $l$  is indeed the equilibrium distribution. This means that its potential on  $S$  is a constant that we may assume to be equal to 1:

$$\iint_S \frac{l(x')}{|x - x'|} d\omega(x') = 1. \quad (15)$$

In order to abbreviate the rest of the argument, and to see the connection with modern spectral theory I shall introduce some notation, that cannot be found by Liouville. Consider the Hilbert space  $L^2(S, l d\omega)$  with the inner product

$$\langle \zeta, \eta \rangle = \iint_S \zeta(x) \eta(x) l(x) d\omega(x). \quad (16)$$

Moreover define the integral operator  $A$  on  $L^2(S, l d\omega)$  by

$$A\zeta = \iint_S \frac{\zeta(x')}{|x - x'|} l(x') d\omega(x'). \quad (17)$$

In this terminology (15) can be written as

$$A1 = 1, \quad (18)$$

i.e. the function 1 is an eigenfunction of  $A$  with eigenvalue 1. In the second step Liouville finds the next eigenfunction and eigenvalue by maximizing  $\langle A\zeta, \zeta \rangle$  under the condition that  $\langle 1, \zeta \rangle = 0$  and  $\|\zeta\|^2 = \text{constant}$ . The maximizing function  $\zeta_1$  has

$$A\zeta_1 = m_1\zeta_1 \quad \text{where } m_1 \leq 1. \quad (19)$$

Generally Liouville's  $n$ 'th step calls for the maximization of  $\langle A\zeta, \zeta \rangle$  in the orthogonal complement of the subspace spanned by the eigenfunctions found in the  $n - 1$  previous steps. The maximizing function  $\zeta_n$  has

$$A\zeta_n = m_n\zeta_n \quad \text{where } m_n \leq m_{n-1} \leq \cdots \leq m_1 \leq 1 \quad (20)$$

It is clear that these eigenfunctions are mutually orthogonal, and in the last lines of the quote Liouville claims that any function on  $S$  can be expanded in a Fourier series of the eigenfunctions.

The above variational technique for finding eigenfunctions is now the standard method for compact operators and is often attributed to Rayleigh (1877). Liouville's note raises many questions:

- 1) What motivated Liouville to develop this method?
- 2) How much of a proof did he give?
- 3) Why did he not publish his results?
- 4) Are Liouville's results correct?
- 5) What are the historical and mathematical connections to later works?

#### 4.3 Motivation

Liouville was motivated mainly by Gauss' great paper on potential theory (1840) in which the following problem was formulated:

**Gauss' Problem.** *Given a function  $\bar{V}$  on a surface  $S$ . Determine a single layer distribution  $\lambda$  on  $S$  whose potential has the value  $\bar{V}$  on  $S$ .*

In 1842, in connection with his (also unpublished) research on stability of rotating masses of fluid Liouville introduced the Lamé functions  $S_{iB}$  of the second kind, in addition to Lame's  $R_{iB}$ ,  $M_{iB}$  and  $N_{iB}$ . The only thing we need to know about these functions is that they are defined in  $\mathbb{R}^3$  in such a way that  $R_{iB}$  and  $S_{iB}$ , for each value of  $i$  and  $B$ , are constant on a certain ellipsoid  $S$ , whereas  $M_{iB}$  and  $N_{iB}$  can be considered as functions on this surface. In 1845 and 1846 Liouville published a series of results about the Lamé functions. The main theorem states that:

$$\iint_S \frac{l(x')M_{iB}(x')N_{iB}(x')}{|x - x'|} d\omega(x') = \frac{4\pi R_{iB}(x)S_{iB}(x)}{2i + 1} M_{iB}(x)N_{iB}(x). \quad (21)$$

Here  $l$  is a function on  $S$  that Liouville later interpreted as the equilibrium charge distribution. Liouville remarked, that if  $\zeta_n$  denotes the product  $M_{iB}N_{iB}$ ,  $\zeta_n$  will satisfy the eigenvalue equation

$$\iint_S \frac{l(x')\zeta_n(x')}{|x - x'|} d\omega(x') = m_n \zeta_n(x). \quad (22)$$

It is now easy to solve Gauss’ problem for the ellipse. Indeed if we write  $\bar{V}$  as

$$\bar{V} = \sum_n A_n \zeta_n \quad (23)$$

(Liouville showed that this could “always” be done), then

$$\lambda = l \sum_n \frac{A_n}{m_n} \zeta_n \quad (24)$$

is the desired charge distribution. The proof is an easy application of (22) and the orthogonality of the  $\zeta_n$ ’s.

At this point Liouville got the great idea of generalizing this method to an arbitrary surface  $S$ , i.e. to find its equilibrium distribution  $l$ , and define the eigenfunctions  $\zeta_n$  by (22), such that any function  $\bar{V}$  can be expanded in a series (23). Then  $\lambda$  defined by (24) is a solution of Gauss’ problem. He sketched this idea briefly in his published paper (1845a) and even emphasized its importance:

After having studied the matter I do not hesitate to regard the functions  $\zeta$  as being of the utmost importance in analysis (Liouville 1845a).

The problem raised in Liouville’s notebooks is the following: How can the  $\zeta_n$ ’s be found? As we saw, the variational method provided the answer. His notebooks even reveal how he arrived at the method: First he used the idea *a posteriori*: Assuming that the  $\zeta_n$ ’s exist with  $1 > m_1 > m_2 > \dots$  and supposing that any function can be expanded as  $\sum A_n \zeta_n$ , it is easy to see that  $m_n$  and  $\zeta_n$  can be found by the variational method. Some 40 pages later he returned to this method and remarked that “considered *a priori* it demonstrates the existence of these functions”.

#### 4.4 Proofs

In my Liouville biography (Lützen 1990 Chap. XV) one can find as much of Liouville’s proofs as I have been able to reconstruct from Liouville’s rather disorganized notes. These proofs are really beautiful up to a certain point and are strikingly similar to the arguments found in a modern book on spectral theory of compact operators. Liouville succeeded in proving Bessel’s inequality, Parseval’s equality and showed that if a function is orthogonal to all the  $\zeta_n$ ’s it must vanish identically. From this he concluded that if the Fourier series of a function  $Q$  converges, it has the sum  $Q$ . However, since he did not have the Lebesgue integral at his disposal, and since he was interested in pointwise or uniform convergence rather than  $L^2$  convergence, he could not prove that the

Fourier series converges. Thus the situation was the oposite of the situation in Sturm-Liouville theory, where he could prove the convergence of the Fourier series but could not show that it had the desired sum.

#### 4.5 Why Were Liouville's Results Never Published?

Liouville developed several interesting theories and results that he never published. In general the main reason for not publishing was lack of time, but in the case we consider here, there were also inner mathematical reasons. Indeed there were holes in Liouville's argument. First, he does not seem to have been able to prove that the decreasing series of eigenvalues  $m_n$  tend to zero, and this is of great importance for the whole theory. Liouville may not have been aware of this shortcoming but toward the end of his notes he began to sense another problem, in connection with the variational technique for finding  $l$  and the eigenfunctions  $\zeta_n$ . Indeed he discovered that though the upper (lower) bound exists the max (min) may not always be attained; instead there might be what he called a "tendance indéfinie vers le but".

It is very remarkable that Liouville here pointed to the weakness of the naive version of the Dirichlet-principle, and similar principles, about 25 years before Weierstrass presented his famous counterexample (1870–1895). Of course, once the variational method was called into question the existence of the  $\zeta_n$ 's was not secured, and this may explain why Liouville did not publish his ideas.

#### 4.6 Are Liouville's Results Correct?

This question is urgent, once we have realized that Liouville's arguments are not entirely satisfactory. In modernized language it is enough to show that Liouville's operator  $A$  is compact on  $L^2(S, l d\omega)$  and, indeed, Christian Berg has succeeded in proving this theorem under rather weak assumptions on the surface  $S$  (Berg and Lützen 1990).

#### 4.7 Connections to Later Works

In fact it was not Rayleigh but Heinrich Weber (1869) who first published a variational method for finding eigenfunctions. He applied the method to the two dimensional Laplacian, but as Liouville, he did not show that the infima were attained by minimizing functions. It was left to Poincaré (1894) to establish Weber's results rigorously. In connection with Liouville's notes a later paper by Poincaré (1896) is even more interesting. Here Poincaré introduced what he called fundamental functions of a surface  $S$ . They were defined by a variational method, and are in fact identical to Liouville's  $\zeta_n$ . Poincaré was unable to prove rigorously that these functions exist, and so he only used them as a heuristic tool. However his student Le Roy (1898) succeeded in giving a very long proof of their existence.

Weber's, Rayleigh's, and Poincaré's work, together with works by Kirchhoff, Klein and others from arround 1880 mark the first continuation of Sturm-Liouville theory, in the sence that they dealt with differential equations. Integral operators (or equations) implicitly arose out of the Beer-Neumann method for

solving the Dirichlet problem. They were studied by Volterra in the 1890s but their spectral theory was not studied until 1900 in an elegant paper by Fredholm.

Fredholm’s ideas were carried much further in Hilbert’s monumental work (1904–1910), that in a sense marks the conclusion of the development we have considered. Hilbert introduced the so-called “vollstatige” operators, and showed that their eigenvalues can be found by the variational technique used by Liouville. Moreover he joined the two types of spectral theory that Liouville had studied. Indeed he showed how one can use a Green’s function to reduce problems in Sturm-Liouville theory to (easier) problems of integral operators. He even connected both these theories to spectral theory of discrete systems, by showing how the spectral theorem of quadratic forms in the limit would lead to the spectral theorem of integral operators, or more generally of quadratic forms in infinitely many variables. This idea had already been suggested by Poincare but Hilbert was the first to carry out the limit procedure rigorously. In the case of infinite quadratic forms he made an ingenious use of the Stieltjes integral in order to deal with the continuous spectrum.

About the same time, the structural movement made another unification of the three branches of spectral theory possible. It became clear that they all dealt with one and the same question for different *operators* between different *spaces*. Liouville clearly did not have such a unifying concept, but there is no doubt that he was aware of the formal similarity between his work on Sturm-Liouville theory and his unpublished work on integral operators. Indeed many of the theorems are the same and their proofs often go along the same lines.

## 5. Conclusion

One may wonder if the development of spectral theory would have been speeded up if Liouville had published his research on his integral operator in 1846. This is not at all certain. In fact, Sturm-Liouville theory was not developed further until half a century after Sturm and Liouville had published their work, and Liouville’s small published note on symmetric integral operators went completely unnoticed. How can this lack of interest be explained?

The reason is probably the very qualitative nature of the results. I have stressed this in connection with Sturm-Liouville theory and it is equally true of the results concerning Liouville’s operator  $A$ . In fact, except for simple surfaces, even the equilibrium distribution  $l$  is impossible to determine. Most mathematicians in the middle of the 19th century were not excited by existence theorems, unless they could find the objects. Therefore they did not continue Sturm-Liouville theory and would probably have neglected Liouville’s spectral theory of his integral operator even if it had been published. Only at the end of the 19th century, the mathematical community was ripe for such questions.

This shows Liouville’s farsightedness.

Liouville is usually mentioned in histories of spectral theory in connection with his work on Sturm-Liouville theory. I hope this paper has demonstrated, that he was an even greater pioneer in this field than it is usually acknowledged.

## References

- Abel, N.H. (1823): Oplösning af et Par Opgaver ved Hjelp af bestemte Integraler. Magazin for Naturvidenskaberne, ser. 1, I (1923), 11–27. [French translation in *Oeuvres*]
- Bang, T., Fuglede, B. (1990): No two quotients of Normalized Binomial Mid-coefficients are equal. *J. Number Theory* **35**, 345–349
- Berg, C. (1990): Integrals involving Gegenbauer and Hermite polynomials on the imaginary axis. *Equationes Math.* (to appear)
- Berg, C., Fuglede, B. (1990): Liouville's operator for a disc in space. *Manuscr. math.* **67**, 165–185
- Berg, C., Lützen, J. (1990): J. Liouville's Unpublished Work on an Integral Operator in Potential Theory. A Historical and Mathematical Analysis. *Expositiones Mathematicae* **8**, 97–136
- Cauchy, A.L. (1824–1881): *Équations Différentielles Ordinaires*. Etudes Vivantes, Paris. New York 1981. Fragment of lecture notes for the second year at the École Polytechnique. Introduction by C. Gilain
- Dieudonné, J. (1981): History of Functional Analysis. North-Holland, Amsterdam 1981
- Dirichlet, P.G. Lejeune (1829): Sur la convergence des séries trigonométriques qui servent à représenter une fonction arbitraire entre des limites données. *J. Reine Angew. Math.* **4**, 157–169; *Werke* **1**, 117–132
- Fourier, J. (1822): *Théorie analytique de la chaleur*. Paris 1922
- Fredholm, I. (1900): Sur une nouvelle méthode pour la résolution du problème de Dirichlet. *Öfversigt af Kungliga Svenska Vetenskabs-Akademiens Förhandlanger*, Stockholm **57**, 39–46; *Oeuvres Complètes*, 61–68
- Gauss, C.F. (1840): Allgemeine Lehrsätze in Beziehung auf die im verkehrten Verhältnisse des Quadrats der Entfernung wirkenden Anziehungs- und Abstossungs-Kräfte. Resultate aus den Beobachtungen des magnetischen Vereins im Jahre 1839 (4) Leipzig 1840; *Werke* **5**, 197–242; Ostwald's Klassiker Nr. 2, Leipzig 1889. [French transl. in *J. Math. Pures Appl.* **7** (1842) 273–324]
- Hawkins, T. (1975): Cauchy and the Spectral Theory of Matrices. *Historia Mathematica* **2**, 1–29
- Hilbert, D. (1904–1910): Grundzüge einer allgemeinen Theorie der linearen Integralgleichungen. *Nachrichten Königl. Ges. Wiss. Göttingen Math.-Phys. Kl.* (1904) 49–91; 213–259; (1905) 307–338; (1906) 157–227; 439–480; (1910) 355–417; 595–618. Collected as monograph Leipzig 1912. References refer to monograph
- Le Roy E. (1898): Sur l'intégration de l'équation de la chaleur (2<sup>e</sup> partie). *Ann. Ec. Norm. Sup.* (3) **15**, 9–178
- Liouville, J. (1830): Mémoire sur la théorie analytique de la chaleur. *Ann. Math. Pures Appl.* **21** (1830–1831) 131–181; Summary by Sturm in *Bull. Sci. Math. Phys. et Chim.*
- Liouville, J. (1832): Sur quelques questions de géometrie et de mécanique, et sur un nouveau genre de calcul pour résoudre ces questions. *J. Ec. Polyt.* **13** (21. cahier) 1–69
- Liouville, J. (1836): Mémoire sur le développement des fonctions ou parties de fonctions en séries, dont les divers termes sont assujettis à satisfaire à une même équation différentielle du second ordre, contenant un paramètre variable. *J. Math. Pures Appl.* **1**, 253–265; *C. R. Acad. Sci. Paris* **1** (1835) 418
- Liouville, J. (1837a): Second Mémoire sur le développement des fonctions ou parties de fonctions en séries, dont les divers termes sont assujettis à satisfaire à une même équation différentielle du second ordre, contenant un paramètre variable. *J. Math. Pures Appl.* **2**, 16–35

- Liouville, J. (1837b): Troisième Mémoire sur le développement des fonctions ou parties de fonctions en séries, dont les divers termes sont assujettis à satisfaire à une même équation différentielle du second ordre, contenant un paramètre variable. *J. Math. Pures Appl.* **2**, 418–437; *C. R. Acad. Sci. Paris*, **5**, 205–207
- Liouville, J. (1838): Premier mémoire sur la théorie des équations différentielles linéaires, et sur le développement des fonctions en séries. *J. Math. Pures Appl.* **3**, 561–614; *C. R. Acad. Sci. Paris*, **7**, 1112–1116
- Liouville, J. (1845a): Sur diverses questions d'analyse et de physique mathématique. *J. Math. Pures Appl.* **10**, 222–228
- Liouville, J. (1845b): Sur une propriété générale d'une classe de fonctions. *J. Math. Pures Appl.* **10**, 327–328
- Liouville, J. (1846): Lettres sur diverses questions d'analyse et de physique mathématique, concernant l'ellipsoïde, adressées à M.P.H. Blanchet – Première Lettre. *J. Math. Pures Appl.* **11**, 217–236. Deuxième Lettre. *J. Math. Pures Appl.* **11**, 261–290
- Liouville, J. (Ms 3615–3640) The Liouville Nachlass at the Institut de France
- Lützen, J. (1982): Joseph Liouville's Contribution to the Theory of Integral Equations. *Historia Mathematica* **9**, 373–391
- Lützen, J. (1984): Sturm and Liouville's Work on Ordinary Linear Differential Equations. The Emergence of Sturm-Liouville Theory. *Arch. Hist. Exact. Sci.* **29**, 309–376f
- Lützen, J. (1987): The Solution of Partial Differential Equations by Separation of Variables: A Historical Survey. *Studies in Mathematics* **26** (ed. E. Phillips) 242–277. Math. Ass. of America
- Lützen, J. (1990): Joseph Liouville (1809–1882). Master of Pure and Applied Mathematics. Springer, New York
- Poincaré, H. (1894): Sur les équations de la physique mathématique. *Rend. Circ. Mat. Palermo* **8**, 57–155; *Oeuvres* **9**, 123–196
- Poincaré, H. (1896): La méthode de Neumann et le problème de Dirichlet. *Acta Math.* **20**, 59–142; *Oeuvres* **9**, 202–272
- Poisson, S.D. (1826): Note sur les racines des équations transcendantes. *Bull. Soc. Philomatique* 145–148
- Rayleigh, Lord (1877): *The Theory of Sound*. London
- Sturm, C. (1836a): Mémoire sur les Équations différentielles linéaires du second ordre. *J. Math. Pures Appl.* **1**, 106–186
- Sturm, C. (1836b): Mémoire sur une classe d'Équations à différences partielles. *J. Math. Pures Appl.* **1**, 373–444
- Weber, H. (1869): Ueber die Integration der partiellen Differentialgleichung  $\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + k^2 u = 0$ . *Math. Ann.* **1**, 1–36
- Weierstrass, K. (1870–1895): Über das sogenannte Dirichletsche Princip. (Read in 1870); *Werke*, Berlin 1895, vol. 2, pp. 49–54



# Mathematics as Metaphor

*Yuri Ivanovich Manin\**

Steklov Mathematical Institute, 42 Vavilova, 117966 GSP-1, Moscow, USSR

あ  
か  
あ  
か  
や  
あ  
か  
あ  
か  
あ  
か  
あ  
か  
あ  
か  
あ  
か  
や  
あ  
か  
あ  
か  
月  
や

*Ordre. [...] Je sais un peu ce que c'est et combien peu de gens l'entendent. Nulle science humaine ne le peut garder. Saint Thomas ne l'a pas gardé. La mathématique le garde, mais elle est inutile en sa profondeur.*

Pascal, Pensées

Myōo (1173–1232)

## Introduction

When H. Poincaré first published in 1902 his book *La Science et l'hypothèse*, it became a bestseller. The first chapter of this book was devoted to the nature of the mathematical reasoning. Poincaré discussed an old philosophical controversy whether the mathematical knowledge could be reduced to long chains of tautological transformations of some basic (“synthetic”) truths or it contained something more. He argued that the creative power of mathematics was due to a free choice of the initial hypotheses-definitions which were later on constrained by the comparison of deductions with the observable world.

The society of our days seems to be much less interested in the philosophical subtleties than Poincaré's contemporaries. I do not want to say that science itself became less popular. Such books as S. Weinberg's *The first three minutes* and S.W. Hawking's *A brief history of time* are sold by hundreds of thousands and favorably reviewed in widely distributed newspapers. What has changed is the general mood. The paradoxality of the new physical theories is perceived less dramatically and more pragmatically. (We can note that the perception of visual arts evolved in much the same way: if the first exhibitions of Impressionists were a kind of spiritual revolution, each new wave of the post-war avant-garde immediately acquired family traits of academism).

---

\* Delivered by Barry Mazur.

In this atmosphere, the heated discussions of the bygone days on the foundational crisis of mathematics and the nature of infinity seem almost irrelevant and certainly inappropriate. The audience responds much livelier to the opinions about school education or a new generation of computers.

This is why I have decided to present at this section an unpretentious essay in which our science is considered as a specialized dialect of the natural language, and its functioning as a special case of speech. This implies certain suggestions about the high school and University training.

## Metaphorism

The word “metaphor” is used here in a non-technical sense, which is best rendered by the following quotations from James P. Carse’s book *Finite and Infinite Games*:

“Metaphor is the joining of like to unlike such that one can never become the other.”

“At its root all language has the character of metaphor, because no matter what it intends to do, it remains language, and remains absolutely unlike whatever it is about”.

“The unspeakability of nature is the very possibility of language”.

Considering mathematics as a metaphor, I want to stress that the interpretation of the mathematical knowledge is a highly creative act. In a way, mathematics is a novel about Nature and Humankind. One cannot tell precisely what mathematics teaches us, in much the same way as one cannot tell what exactly we are taught by “War and Peace”. The teaching itself is submerged in the act of re-thinking this teaching.

This opinion seemingly disagrees with the time-honored tradition of applied mathematics in scientific and technological calculations.

In fact, I want only to restore a certain balance between the technological and the humanitarian sides of mathematics.

## Two Examples

Let me try to illustrate the metaphoric potential of mathematics by discussing two disjoint subjects: the Kolmogorov complexity and the “Dictator Theorem” due to K. Arrow.

i) Kolmogorov’s complexity of a natural number  $N$  is the length of a shortest program  $P$  that can generate  $N$ , or the length of a shortest code of  $N$ . A reader should imagine a way of coding integers which is a partial recursive function  $f(P)$  taking natural values. Kolmogorov’s theorem states that among all such functions there exist the most economical ones in the following sense: if  $C_f(N)$  is the minimal value of  $P$  such that  $f(P) = N$ , then  $C_f(N) \leq \text{const} \cdot C_g(N)$ , where const. depends only on  $f, g$  but not  $N$ .

Since  $P$  can be reconstructed from its binary notation, the length  $K_f(N)$  of the shortest program generating  $N$  is bounded by  $\log_2 C_f(N)$ . This function, or rather

the class of all such functions defined up to a bounded summand, is the Kolmogorov complexity.

First of all,  $K(N) \leq \log N + \text{const}$ . Of course, this conforms nicely with the historical successes of the positional notation systems which provided us with the number generating programs of logarithmic length. However, there are arbitrary large integers whose Kolmogorov's complexity is much smaller than the length of their notation, e.g.,  $K(10^N) \leq K(N) + \text{const}$ . In general, when we use large numbers at all, we seemingly use only those which have a relatively small Kolmogorov's complexity. Even decimal decompositions of  $\pi$  which are, probably, the longest well-defined numbers ever produced by mathematicians, are Kolmogorov simple, because  $K([10^N\pi]) \leq \log N + \text{const}$ . In general, small Kolmogorov complexity = high degree of organization.

On the other hand, almost all integers  $N$  have the complexity close to  $\log N$ . For example, if  $f(P) = N$  for an optimal  $f$ , then  $K(P)$  is equivalent to  $\log P$ . Such integers have many remarkable properties which we usually connect with "randomness".

Second, Kolmogorov's complexity can be easily defined for discrete objects which are not numbers, for example, Russian or English texts. Therefore, "War and Peace" has a pretty well defined measure of its complexity; the indeterminacy is connected with the choice of an optimal coding and seems to be pretty small if one chooses one of the small number of reasonable codings.

From this viewpoint, is "War and Peace" a highly organized or an almost random combinatorial object?

Third, Kolmogorov's complexity is a non-computable function. More precisely, if  $f$  is optimal, there is no recursive function  $G(N)$  which would differ from  $C_f(N)$  by  $\exp(O(1))$ . One can only bound complexity by computable functions.

I feel that Kolmogorov's complexity is a notion that is very essential to keep in mind in any discussion of the nature of human knowledge.

As long as the content of our knowledge is expressed symbolically (verbally, digitally, ...) there are physical restrictions on the volume of information that can be kept and handled. We always rely upon various methods of information compressing. Kolmogorov complexity puts absolute restrictions on the efficiency of such a compression. When we speak, say, of physical laws, expressed by the equations of motion, we mean that a precise description of the behaviour of a physical system can be obtained by translating these laws into a computer program. But the complexity of laws we can discover and use is clearly bounded. Can we be sure that there are no laws of arbitrary high complexity, even governing the "elementary" systems?

At this point, our discussion becomes totally un-mathematical, and before a mathematically-minded audience I must stop here. But such is the fate of any metaphor.

ii) Arrow's Dictator Theorem was discovered around 1950. Mathematically, it is a combinatorial statement describing certain functions with values in binary relations. Intuitively, it is a formalized discussion of the problem of Social Choice. Suppose that a lawmaker has to establish a law which governs the processing of

individual wills of voters into a collective decision. If the problem is to choose one of the two alternatives, the standard solution is do it by the majority of votes. However, usually there are more than two alternatives (imagine the funds allocation problems), and voters may be asked to order them according to their preferences. What should be the algorithm extracting the collective preference from any set of individual preferences? Arrow considered algorithms satisfying some natural and democratic axioms (e.g., when everybody prefers *A* to *B*, the society prefers *A* to *B*). Nevertheless, he discovered that when there are more than two alternatives, the only way to achieve a solution is to nominate a member of the society ("the Dictator") and to equate his personal preference order to the social one. (Actually, this is one of the versions of Arrow's theorem discovered later. Also, it refers to the case of a finite society; in the infinite case, the social decisions can be made by ultrafilters, appropriately called "the ruling hierarchies".)

In a way, this theorem illustrates the content of Jean-Jacque Rousseau's idea of a *Contrat Social*.

The fundamental inconsistency of the image of the ideal democratic choice can be illustrated by the following story referring to three voters and three alternatives. It is the story of three knights errant at the cross-roads with a stone before them. The inscription on the stone prophesies only losses: who goes to the left will lose his sword; who goes to the right will lose his horse; who goes straight will lose his head. The knights dismount and start taking council. In a Russian version of this story, the knights have names and personalities: the youngest and ardent Alyosha Popovich, the eldest and wisest Dobrynya Nikitich, and the slow peasant Ilya Muromets. So Alyosha values sword more than horse, and horse more than his head; Dobrynya values most his head, then his sword, then his horse; and Ilya prefers his horse to head to sword.

A reader will note that the three individual preference orders constitute one and the same cyclic order on the set of alternatives. As a result, one can decide by majority the choice between any two of the alternatives, but the union of these decisions will be inconsistent: the democratic procedure cannot provide us with a well-ordered list. The knights sigh and delegate the decision-making power to Dobrynya.

Does the Arrow theorem tell us something that we did not know beforehand? Yes, I think it does if we are ready to discuss it seriously, that is, to look closely at the combinatorial proof, to imagine the possible real life content of various assumptions and elementary logical steps made on the way, in general, to enhance our imprecise imagination by the rigid logic of a mathematical reasoning. We can understand better, for example, some tricks of policy-making and some pitfalls the society can leap whole-heartedly into (like accepting without questioning a list of alternatives imposed by a ruling hierarchy, although precisely the compilation of this list can be the central issue of the social decision making).

At this stage, we come to the main topic of our discussion: what distinguishes a mathematical discourse from a natural language discourse, why the Pascalian "ordre" came to reign over our specialized symbolic activities, and is it truly so "useless in its profundity"?

## Language and Mathematics

A very interesting chapter of the interaction between mathematics and humanities started about thirty years ago when the first serious attempts of automatic translation were made. These first attempts were a painful failure, at least so for many optimists who believed that in this domain, there are no fundamental obstacles, and it remains only to overcome technical difficulties connected with the sheer amount of information to be processed. In other words, they took for granted that the translation is in principle performed by a not very complex algorithm which only must be made explicit and then implemented as a computer program.

This assumption is a nice example of a mathematical metaphor (actually, a specialization of the general “computer metaphor” used in the brain sciences).

This metaphor proved to be extremely fruitful for the theoretical linguistics in general because it forced linguists to start describing vocabulary, semantics, accidence, and syntax of human languages with unprecedented degree of explicitness and completeness. Some totally new notions and tools were discovered thanks to this program.

However, the successes of the automatic translation itself were (and still are) scanty. It became clear that written human speech is an extremely bad input data for any algorithmic processing planned as translation or even as a logical deduction. (I add this proviso because there is nothing special in human speech considered as a material for, say, statistical studies).

This fact can be considered as a universal property of human languages, and it deserves some attention. One must first of all reject as too naive a usual explanation that the universe of meanings of a human language is too vast and poorly structured to admit a well organized metalanguage describing this world. The point is that even if we severely restrict this universe to the subset of arithmetic of small integer quantities, we shall still have to face the same difficulty. In fact, this difficulty was a decisive reason for the crystallization of the whole system of arithmetical notation and the basic algorithms of calculation, and later on of symbolic algebra. Even the vocabulary of elementary arithmetic in human language is basically archaic: the finite natural series of primitive societies “one, two, three, indefinitely many” is reproduced in the exponential scale in our “thousand, million, billion, zillion”. The expressions for relatively small numbers like “1989” are actually names of the decimal notation and not of the numbers themselves.

The advantage of F. Viète’s algebra over the semi-verbal algebra of Diophantus was due not to the fact that it could express new meanings but to the incomparably greater susceptibility to the algorithmic processing (“identical transformations” of our high school algebra).

The rupture of the intuitive and emotional ties between a text and its producer/user so characteristic of the language of science was compensated by the new computational automatisms. In their (albeit restricted) domain they proved to be infinitely more efficient than the traditional Platonian and Aristotelian culture of everyday language discourse. Why then our scientific papers are still written as a disorganized mixture of words and formulas? Partly because we still need those emotional ties; partly because some meanings (like human values) are best rendered

in human language. But even as a medium of scientific speech, human language has some inherent advantages: appealing to the spatial and qualitative imagination, it helps to understand “structurally stable” properties like the number of free parameters (dimension), existence of extrema, symmetries. To put it bluntly, it makes possible the metaphorical use of science.

## Metaphor and Proof

The views professed here can be considered in relation to the high school and graduate curriculum.

The general mathematical education of the first half of this century was application oriented. It provided the basic minimum for the practical life problems and a smooth transition to the study of engineering and scientific calculations at the college level. The break of this curriculum with the activity of professional mathematicians became more and more pronounced. As is well known, this brought the reaction in the form of NewMath in USA and similar programs in other countries. These programs introduced into the high school mathematics the notions and principles borrowed from professionals: set theory, axiomatic presentation of proofs, strict culture of definitions.

NewMath became widely accepted but its expansion was accompanied by the protesting voices which in the 70s and 80s merged into a loud chorus. The critics disagreed with the basic arguments of the NewMath proponents. Leaving aside the objections based upon the data from cognitive sciences and learning psychology I shall only recall those concerning the general evaluation of the role of the proof in mathematics.

The one pole is represented by the well known statement due to Nicolas Bourbaki: “Depuis les Grécs, qui dit Mathématiques, dit démonstration”. According to this perception, the rigorous proof was made a matter of principle in the NewMath programs. It was argued that: a) a proof helps to understand a mathematical fact; b) a rigorous proof is the most essential component of the modern professional mathematics; c) mathematics possesses the universally recognized criteria of rigour.

These views were extensively criticised, e.g. by Gila Hanna in the book “Rigorous Proof in Mathematics Education”, OISE Press, Ontario 1983. In particular, Gila Hanna pointed out that the mathematicians are far from unanimous in accepting the criteria of rigour (referring to quarrels between logicians, formalists and intuitionists), and that working mathematicians constantly break all rules in the book.

In my opinion, this is irrelevant.

What is relevant, is the imbalance between various basic values which is produced by the emphasis on proof. Proof itself is a derivate of the notion of “truth”. There are a lot of values besides truth, among them “activities”, “beauty” and “understanding”, which are essential in the high school teaching and later. Neglecting precisely these values, a teacher (or a university professor) tragically fails. Unfortunately, this also is not universally recognized. A sociological analysis of the controversies around the Catastroph Theory of René Thom shows that exactly the

shift of orientation from the formal truth to understanding provoked such a sharp criticism. But of course, the Catastroph Theory is one of the developed mathematical metaphors and should only be judged as such.

Pedagogically, a proof is just one of the genres of a mathematical text. There are many different genres: a calculation, a commented sketch, a computer program, a description of an algorithmic language, or such a neglected kind as a discussion of the connections between a formal definition and intuitive notions. Every genre has its own laws, in particular, laws of rigour, which only are not codified because they were not payed a special attention.

A central problem of a teacher is to demonstrate at the restricted area of his or her course the variety of types of mathematical activities and underlying value orientations. Of course, this variety is hierarchically organized. The goals may vary from achieving an elementary arithmetical and logical literacy to programming skills, and from the simplest everyday problems to the principles of modern scientific thinking. In the spectrum of these goals, the emphasis on the norms of "rigorous proof" can safely occupy a peripheral position.

But having said all this, I must stress that my argumentation by no means undermines the ideal of a rigorous mathematical reasoning. This ideal is a fundamental constituting principle of mathematics, and in this sense Bourbaki is certainly right. Having no external object of study, being based on a consensus of a restricted circle of devotees, the mathematics could not develop without the permanent control of rigid rules of game. Applicability of mathematics in the strict sense of this word (like its indispensability in the Apollo project) is due to our ability to control series of symbolic manipulations of fantastic length.

The existence of this ideal is far more essential than its unattainability. The freedom of mathematics (G. Cantor) can only develop in the limits of iron necessity. The hardware of modern computers is an incarnation of this necessity.

Metaphor helps a human being to breathe in this rarefied atmosphere of Gods.



# Teaching Mathematics to Students Not Majoring in Mathematics

– Present Situation and Future Prospects –

*Haruo Murakami*

Department of Applied Mathematics, Faculty of Engineering, and  
Division of Intelligence Science, Graduate School of Science and Technology  
Kobe University, Kobe, 657 Japan

The present era, described variously as a Highly Industrialised Society, the Information Era or even a Postindustrial Society, places greater emphasis than before on the importance of teaching mathematics. Teaching mathematics to students who select it as their major subject is of course a problem of vital importance, but compared with all the students studying mathematics, those actually majoring in it would be a drop in the ocean; the problem of teaching mathematics to those majoring in other subjects is a much larger problem.

University education has become more available and more widely spread. In Japan, for instance, there are three times as many university students as there were middle-school pupils before the 2nd World War, and the present university teaching staff is five times more numerous than pre-war middle-school teachers. Universities have rid themselves of their ivory tower image and have adapted to the need for practical sciences. At present, there is probably more demand for this, a more practically oriented area.

Most of those students who require mathematical education at university level specialise in engineering, the rest in natural sciences such as physics, chemistry, and biology, as well as those who specialise in economics, business studies, or some areas of medicine, etc. In fact, in Japan, there are only 3,200 students specialising in pure mathematics, whereas almost all the students in engineering, who number well over 65,000, receive lectures in mathematics.

## 1. *What kind of mathematics and mathematical teaching is required?*

Mathematics is indispensable for the study of engineering, no matter what the field of engineering and no matter what the area of technology the student selects for his or her future professional skill. Students of engineering cannot therefore avoid learning mathematics. At present, although the extent varies according to the field, all students of engineering get some mathematical tuition. What passes for mathematics here is mathematical skills, i.e. mathematical knowledge and the ability to apply this knowledge to engineering, ability to calculate and compute, etc. The subjects normally taught to them are: calculus, linear algebra, theory of complex functions, vector analysis, special functions such as Bessel functions,

differential equations, probability and statistics, etc. They would be quite unable to come to grips with their field of engineering without knowing these subjects. It is therefore unlikely that the need for teaching these subjects should decrease in the future.

There are naturally some differences depending on the specialisation. For example, those who enter fields where fluid dynamics is important, such as civil and mechanical engineering, may find vector analysis important. But for those who aim at signal theory, it would be essential to study abstract algebra and in particular Galois fields. There are some areas of engineering where topology and graph theory may be more important, or it may be essential to study Boolean algebra. How would one satisfy such diverse demands of engineering students?

One might say that except for engineering, a basic knowledge of mathematics would suffice for any discipline. For an ordinary student in economics, elementary calculus and linear algebra used to be sufficient. However, econometrics, which is increasingly becoming an essential part of economics, demands some knowledge of probability and statistics, and further, depending on one's specialty, a student in economics may have to know such subjects as fixed point theorems and the theory of differential equations. Similarly, the students of physics, chemistry and many other fields all require some degree of specialised knowledge in mathematics.

There are also some subjects in mathematics which seem irrelevant or unimportant to one's own field at first, and yet turn out to be important or even essential later. Since it is impossible to teach all potential areas of mathematics covering many different subjects, we ought to lay a certain amount of groundwork in order to equip students with the ability to handle new subjects as they arise. For example, there was a time when Boolean algebra was not considered a part of applied mathematics. Now, however, Boolean algebra is even more common in applied than in pure mathematics. We shall probably see this sort of thing happening in the future too. We have no idea which techniques will be used and no idea of what to provide in the way of basic education for these new techniques. But to go back to the above example, it is not difficult to deal with Boolean algebra as long as one has some grounding in abstract algebra. This shows how important it is to provide training in the basics of mathematics as a whole. In the information age, beset as it is by violent changes, it is difficult, and even dangerous, to make prior judgments and predictions about the sort of mathematics that will be needed. Surely our aim should be to give our students the ability to learn new topics, and, to enable them to confront new things without fear, although this is easier said than done.

*2. Should one confine oneself to teaching students only the knowledge and skills immediately required for their study, or should they be taught some different aspects of mathematics as well? If so, what should these be?*

If we consider the future development of applied mathematical techniques on the assumption that computers will be more widely used, then we should probably place the emphasis on *intelligence* rather than *skill*, on *thought* rather than *computation*, and on *concepts* rather than *formulae*.

Applied mathematicians of the old school were concerned solely with the memorization and use of formulae. They were, so to speak, *walking reference books* of mathematical formulae. Nowadays, however, many more scientific fields have a very high mathematical content, and mathematics is now applied to a much wider range of subjects. Thanks to modern computer techniques, computational skills will be replaced by the use of computers. Take the example of teaching techniques of integration. The basics such as integration by parts and integration by substitution must be taught as a part of differential and integral calculus, but there is no need to go into finer technical points, or give the students hours of computational practice. It is more efficient to train them to use computer algebra systems or formula manipulation systems such as REDUCE or MATHEMATICA, just as students in the past were trained to use formula books. In the same way, we can delegate the specialised techniques of solving differential equations to computer formula manipulation systems.

In the past, engineering departments tended to view mathematics as something required for physics. The mathematics used in physics generally deals with phenomena which are quantified from the start, and in many cases, it is possible to manage with just classical mathematics or statistics. In recent years, however, even in engineering, there has been a rapid expansion in subjects dealing with social phenomena, such as environmental engineering, urban engineering, and management engineering, most of which deal with unquantifiable things. This has led to a greater need for abstract mathematics. Taking this into account, the emphasis of mathematical education from now on should be placed on training students to become able to grasp abstraction and logical thought.

Many engineers retain the old-fashioned “walking reference book” or “cook-book-method” approach to mathematics. This tends to filter through to engineering students who then attempt to try to find applicable formula from existing lists and apply the result of the computation to an engineering problem. I would like to stress the enormous importance of logical thought through the teaching of mathematics. We should highlight mathematical way of thinking, and concentrate on basic mathematical concepts. Thus, when teaching integration, we should stress the *meaning* of the integral more than how to *evaluate* a given integral. Even when explaining a theorem, we should keep the proof to a minimum. It is more important to stress premises and offer counter-examples where these premises do not hold true.

### 3. How do computers change the teaching of mathematics?

In the face of changing times and especially with the advent of the Information Era, is it still correct to teach the same things as we have done so far, or should we reconsider the contents of our syllabus? How can we introduce the computer more effectively in education? Which new ways of teaching, making full use of computers, can be devised?

The contents of the mathematics to be taught in the computer era must, of necessity, have changed. Discrete mathematics and finite mathematics have

become more important than before. Although more syllabus time should be devoted to these areas, this does not mean that traditional continuous mathematics should lose importance. Moreover, to gain deeper understanding, one also needs these skills for concrete computation and calculation. What exactly are these skills, then? As the expertise required for modern warfare are different from that used during the age of bows and arrows, so the computational skills required in the age of computers are vastly different from those required in the days of pencils and paper. The skill of handling computers has already become a part of a mathematician's skill. There may be other new kinds of skills required in the future.

Originally, mathematics possessed an experimental side. Euler speaks of pure mathematics thus: "The properties of numbers which everybody knows have been discovered through ordinary observation. They were discovered long before their correctness could be verified by proof. It is by observation that more and more properties are discovered, and afterwards people spend all their efforts on proving them." The emergence of the computer expands the possibilities of observation and experimentation in mathematics.

Since it is easy to draw graphs of functions and the solution curves of differential equations on a personal computer, this can be used to enhance learning and research. For instance, it is now easy to draw the solution graphs of a differential equation for as many cases as one wishes, extract those properties which are in common, and make a conjecture on the nature of the solutions. One can then attempt to prove the conjecture, but before reaching this stage, the skill of handling computer graphics is required. There are still some pure mathematicians who have an aversion to computers, but surely this aversion is merely backwardness.

Once upon a time, man used only his own strength, but he then learned to use the strength of domestic animals such as horses to do his labour. Later, engines were developed which were hundreds or thousands of times more powerful than the animals. These engines now routinely do work that would be almost impossible by manpower alone. Life without these machines is virtually unthinkable.

In the same way, man has also learned to use technology to perform his intellectual work. Computers have been developed which can process information hundreds or thousands of times faster than the human brain, and can routinely do work that would be almost impossible by brainpower alone. There are still many areas where computers cannot replace the power of human thought, but they are now becoming essential tools, and it would be difficult to contemplate life without them. Mathematics research and teaching must adapt accordingly in order to benefit from the progress of computers.

At the very least, we could install a microcomputer in the classroom with a big screen to demonstrate effectively the locus of functions, solution curves of differential equations and so on. We could then ask students to do their assignments either by using their own personal computers or by using computing facilities provided by the university. As I mentioned above, a computer algebra system could be employed as a very powerful tool in the mathematics education.

*4. Which subjects should be taught? One example*

One topic proposed as a practical example of mathematical over-abstraction was that old chestnut, the existence theorem of solutions to differential equations. This can be taught with the use of Picard's successive approximation. Of course, instead of carrying out a rigorous proof of the Picard successive approximation, one could just settle it by arm waving. There may be people who think it unnecessary to teach the existence theorem. I believe that it would be better to teach the theorem. One reason is that successive approximation is a concrete way of showing the existence of a solution by construction and this is useful in engineering. Another reason is that the mathematics which students have learned at high school does not prepare them for the possibility that there may be no solution to a problem. When they were given a problem to solve, they could always take it for granted that it had a solution. It would be very useful, for several purposes, to give the students a little culture shock by teaching them the existence theorem. However, those teaching first-year university students should be aware of the fact that their students are completely new to mathematics involving things like existence theorem. Also, topics like this should not take up too much time in the syllabus.

However, it may be a little hard on the students if there are too few subjects of computation, like calculation of indefinite integrals for instance, which follow on from the high-school curriculum, even though it should not be covered in great depth. It is by no means easy for students to overcome the gap between high-school mathematics and university mathematics. Should we change high-school mathematics, then? It would be quite inappropriate to cram subjects like existence theorems wholesale into the high-school curriculum. In fact, mathematics is not only an abstract science, it also is a technique. If it is taught under the sole auspices of mathematicians, they might overimpose their taste for the abstract. The useful aspects may be obscured, and worse still, there is the danger that even the deep understanding of mathematics may be lost. This is because there can be no deep understanding without a fair amount of skill in concrete calculation, etc. I was told that some students of Tokyo University, while knowing all sorts of abstract things, were quite unable to expand  $\sin x$  into Taylor series, and that primary school pupils in France who were asked what " $1 + 2 = ?$ " is, answered " $1 + 2 = 2 + 1$ ," adding that the operation of addition is commutative. Such anecdotes are in fact quite embarrassing.

*5. In which direction should the teaching be aimed? Is it appropriate to choose the same method as applied to teaching mathematics specialists, or is a different approach called for altogether?*

Although mathematics for engineering students should not be very different from that for pure mathematics students, one should be aware of the fact that pure mathematics students are those who are already attracted or enchanted by the beauty of mathematics. These students are willing to study mathematics as their major subject, whereas students of other subjects think of mathematics as merely

a tool to use for their own studies. Thus we need to give the students motivation to study. How can this be achieved? We can begin by showing them how one can construct mathematical models and how one can solve them. What I have been doing in my ODE course for engineering students at Kobe University is this. At the first lecture, they are given a mathematical model of population dynamics. The equation given first is

$$\frac{dx}{dt} = ax.$$

As is well known, this simple Malthus model is not very good at all. They are then introduced to the famous Verhulst model:

$$\frac{dx}{dt} = ax \left(1 - \frac{x}{b}\right).$$

By giving this model, I can talk about the asymptotic behaviour of a solution, equilibrium points, and stability of a solution. Then, a small parameter  $c$  is added.

$$\frac{dx}{dt} = ax \left(1 - \frac{x}{b}\right) - c.$$

This parameter  $c$  may correspond to capturing animals in Africa, or to deforestation by human beings in South America. The parameter  $c$  is then increased a little. Here I am in a position to be able to talk about the structural stability. The parameter can then be increased a little more. If the parameter  $c$  is increased still further, then, all of sudden, the structural stability is lost, and students realise that all the solution curves drop away from the top left to the bottom right. I can then explain what environmental capacity means to us and how easily we can destroy our world.

I have also found it very effective to demonstrate this example of population dynamics by installing a micro-computer in the classroom displaying families of logistic curves on the screen.

Concerning the way of teaching, one can say that there are two ways of teaching, namely, the *top down method* in which general theory is taught first and examples are explained by applying the general theory, and the *bottom up method* in which examples are taught first and general theory is introduced by extracting properties common to these examples. I think the most effective way is perhaps to give a few typical examples like the ones above in order to motivate the students and then formulate a general theory by extracting properties common to these examples. After proving theorems, more examples of application should be given to show how powerful the theory is. This combination of bottom up and top down methods would be the most efficient.

*6. What qualifications should the instructors ideally possess in order to offer the most effective instruction? Should the teaching be undertaken by those who have specialised in pure mathematics and have carried out some research, or should it*

*be undertaken by experienced users of mathematics? Or, indeed, is the ideal tutor for this kind of education someone who is primarily a teaching specialist, research being of secondary concern?*

On this matter, Dr. Satsuma of Tokyo University carried out an interesting survey some years ago. According to his report, the students expressed the following views in response to the question: "Who would be an ideal tutor of mathematics to engineering students?" "Those who once worked hard to learn mathematics because they understand how we feel and how much we suffer from an allergy to mathematics;" "Get rid of those who lecture from inside their own world, while ignoring students;" "Mathematics professionals are likely to stay wrapped up in their cocoon, so give us a non-mathematician;" "An expert on teaching mathematics, who understand those who do not understand."

These comments provoke a little awakening. Why should it be all right for middle-school or high-school teachers to be teaching experts, but when it comes to university, those who do no research at all are suddenly able to teach well? Surely, someone who has never done any research cannot fully understand its importance and, therefore, cannot teach it properly. This would surely halve the value of university education. This is why university teachers must be researchers. Even if this researcher were to lack eloquence, his pervading dedication to research is bound to come across and inspire students. Of course, if this researcher also were to strive to "understand those who do not understand," he or she would become quite outstanding. There are those who would leave research to academies and research institutions and have universities solely as organs of education. This, for the above-stated reasons, would not seem appropriate.

*7. It is safe to assume that the importance of teaching mathematics to students of non-mathematics majors will further increase in the future. What sort of policies and methods have we to cope with this increased importance?*

Mathematics itself, to be sure, has undergone many changes. Not only has there been an expansion in the way mathematics provides the basic means of expression for other sciences but, at the same time, it has itself undergone internal changes. Some famous long standing problems, such as the four colour problem or the continuum hypothesis, have already been solved. At the same time, mathematics has become extremely complex and greatly more diversified. It is now almost impossible to understand a paper on a subject completely unrelated to one's area of study. Small wonder that non-mathematicians regard mathematics as a secret sect of mystics, with some few high priests allowed to carry out secret rituals after long years of specialised study, totally incomprehensible to ordinary mortals. Abstraction has advanced to this degree, particularly since the beginning of this century. A good example of this is the Bourbaki group of French mathematicians in the 1940s, who advanced the rigorousness of expressions. Thus mathematics proceeds to greater refinement, creating a beautifully elegant logical structure – the Temple of Pure Mathematics. This results in weakening the link with the other sciences and engineering, for which mathematics ought to provide the basic

mode of expression. Even the partition walls between mathematics and physics, and also between mathematics and engineering have grown wider.

This trend, however, has changed a little in the last few years. For example, the latest developments in mathematical physics, particularly the advancements concerning nonlinear problems, are an indication of the fact that the two disciplines are again beginning to come closer together. An involvement of algebra and geometry becomes also noticeable, as well as the traditional analysis.

Since mathematics is different from engineering, it is not desirable that only engineers are to be in charge of teaching mathematics at the faculty of engineering. On the other hand, as transpired from the survey mentioned earlier, if somebody who studied mathematics and does research in pure mathematics is put in charge of teaching students of applied sciences, such a tutor has the tendency to become blinkard and teach pure mathematics, as would be taught to students of mathematics. This occasionally confuses the students.

What, then, is the best way of instructing? Those who teach mathematics to students majoring in other subjects, such as engineering, should ideally understand engineering to some extent, so that they may sufficiently comprehend the attitudes and the enthusiasm which students of engineering have towards their subject. At the same time they should have a sound grasp of mathematics as such and should do research in mathematics or in a very closely related subject. Of course, it is almost impossible for one person to know all branches of engineering in detail, to know what is easy and what is difficult, to know what is important and what is trivial, and at the same time have a top-class grasp of pure mathematics. Nonetheless, mathematics should be taught to engineering students by someone who is not entirely ignorant of what engineering is, who has an adequate understanding of the role mathematics plays within engineering, and who knows how one should teach mathematics to them. Unfortunately, it is difficult to see how people of such profile might easily come out of the present educational establishment in any significant numbers.

In consequence, we should consider some effective strategies for producing teachers of that kind. As I mentioned before, both in engineering and in other applied sciences, there has emerged an increasing number of fields which positively utilise mathematics and entirely depend on it. This now covers a very wide area and what is common to this area is that it no longer relies solely on recalling theorems or formulae, but relies more on a thorough understanding of the underlying mathematical structure. This is applied mathematics in the contemporary sense. It lies between pure mathematics and engineering, bridging the gap between the two. This is why applied mathematics has suddenly acquired much more importance. Reflecting this new importance, Japan SIAM (Japan Society of Industrial and Applied Mathematics) was established earlier this year.

I would now like to make a proposal, which I address firstly to the government of Japan. In each faculty of engineering, create a department of applied mathematics or a department of mathematical engineering that would be responsible for the teaching of mathematics in the faculty and would also create a base for research in applied mathematics and mathematical engineering. Since this

department would be involved with teaching mathematics to students of other departments, the staff-student ratio in the department should be appropriately weighted. Those who graduate from such a department are likely to have a foot both in engineering and in mathematics, and should therefore understand both ways of thinking. They would therefore be the ideal graduates for our job. Such people may exceptionally come from computer science or from informatics, but not as a rule. Preferably, they should be inclined more toward mathematics. However, as the role of the computer in those disciplines is quite extensive, the tutors should not suffer from computer allergy.

I should furthermore like to add a few words concerning the staff composition of this department. First, it is likely that a well balanced team could be appointed containing the right number of pure mathematicians in proportion to others. As the number of graduates of the department increased, these graduates may begin occupying a larger part of the staff. At that time there may arise the danger that this would lead to a weakening of the feeling for the mathematical spirit. Unlike other departments within engineering, this department would depend on retaining the mathematical spirit, i.e. mathematics as such. There should always be a reasonable number of pure mathematicians on the staff of these departments if we are to avoid a deterioration. Sufficient attention must be paid to this point.

The situation in other countries is probably similar. It is not a bad idea to establish a faculty of mathematical sciences consisting of a department of pure mathematics, applied mathematics, statistics, and possibly mathematical informatics. Also, it may be thought, for instance, that in developing countries, where the task of developing national resources and industrialisation is pre-eminent, the relevance of applied mathematics far outstrips that of pure mathematics. In these countries, it might be better to begin with application-oriented mathematics and hope that a move towards pure mathematics will take place.

Naturally, we would not be creating this department solely for the purpose of training teachers of mathematics. Most of the graduates from this department would seek employment in industry. It is said that the recent trend in industry is from heavy-thick-long-big to light-thin-short-small. It is impossible to know how long this will continue, but it has resulted in mathematics being directly, rather than indirectly, involved in industry. The mathematics involved is, of course, mainly applied mathematics. Thus the graduates of the departments of applied mathematics or of mathematical engineering that I am proposing here would be the very people now required in industry and would play an important role there. In fact, those countries and industries which do not have such people may well be jeopardizing their future, just as the structure of industries is rapidly upgraded. I should therefore like to emphasise to the industrialists the importance of creating such departments.



# Author Index

I indicates Volume I while II indicates Volume II.

- Alon, Noga 1421-II  
Atiyah, Michael 31-I  
Babai, László 1479-II  
Barbasch, Dan 769-II  
Barlow, Martin T. 1025-II  
Baxter, R. J. 1305-II  
Bedford, Eric 847-II  
Birman, Joan S. 9-I  
Bloch, Spencer 43-I  
Blum, Lenore 1491-II  
Bonahon, Francis 599-I  
Camacho, César 1235-II  
Cameron, Peter J. 1431-II  
Carleson, Lennart 1241-II  
Carlson, Jon F. 317-I  
Chistov, Alexandre L. 1509-II  
Christ, Michael 859-II  
Christodoulou, Demetrios 1113-II  
Coifman, Ronald R. 879-II  
Cook, Stephen A. 55-I  
Coron, Jean-Michel 1123-II  
Cuntz, Joachim 969-II  
Diaconis, Persi 1037-II  
Dolicher, Sergio 1319-II  
Durrett, Richard 1049-II  
Ecalle, Jean P. 1249-II  
Faddeev, Ludwig D. 27-I  
Feigin, Boris L. 71-I  
Feldman, Joel 1335-II  
Floer, Andreas 87-I  
Fukaya, Kenji 491-I  
Furstenberg, Hillel 1057-II  
Gabai, David 609-I  
Ghys, Etienne 501-I  
Gillet, Henri 403-I  
Goldwasser, Shafi 1521-II  
Goodwillie, Thomas G. 621-I  
Gordon, Cameron McA. 631-I  
Grigor'chuk, Rostislav I. 325-I  
Grove, Karsten 511-I  
Günther, Matthias 1137-II  
Harder, Günter 779-II  
Harten, Ami 1549-II  
Hironaka, Heisuke 19-I  
Hofer, Helmut 521-I  
Holmes, Philip 1607-II  
Horiuchi, Annick M. 1639-II  
Huneke, Craig 339-I  
Huxley, Martin N. 413-I  
Igusa, Kiyoshi 643-I  
Ihara, Yasutaka 99-I  
Ikawa, Mitsuru 1145-II  
Il'yashenko, Ju. S. 1259-II  
Ivanov, Alexander A. 1443-II  
Jimbo, Michio 1343-II  
Jones, Lowell E. 653-I  
Jones, Vaughan F. R. 121-I  
Karzanov, Alexander V. 1561-II  
Kashiwara, Masaki 791-II  
Kato, Kazuya 419-I  
Kawamata, Yujiro 699-I  
Kemer, Alexander R. 351-I  
Kollár, János 709-I  
Kolyvagin, Victor Aleksandrovich 429-I  
Kotani, Shinichi 1071-II  
Krasny, Robert 1573-II  
Krichever, Igor 1353-II  
Kronheimer, Peter B. 529-I  
Kupiainen, Antti 1363-II  
Kusuoka, Shigeo 1075-II  
Laumon, Gérard 437-I  
Lazarsfeld, Robert K. 715-I  
Le Cam, Lucien M. 1083-II  
Lebeau, Gilles 1155-II  
Lin, Fang Hua 1165-II  
Lions, Pierre-Louis 1173-II

- Lovász, László 37–I, 139–I  
 Lusztig, George 155–I  
 Lützen, Jesper 1651–II  
 Majda, Andrew J. 175–I  
 Manin, Yuri Ivanovich 3–I, 1665–II  
 Margulis, Grigorii A. 193–I  
 Mathieu, Olivier 799–II  
 Matsuki, Toshihiko 807–II  
 McDuff, Dusa 541–I  
 McMullen, Curt 889–II  
 Melrose, Richard B. 217–I  
 Meyer, Yves F. 1619–II  
 Millson, John J. 549–I  
 Mimura, Masayasu 1627–II  
 Mœglin, Colette 815–II  
 Molchanov, Stanislav A. 1091–II  
 Mori, Masatake 1585–II  
 Mori, Shigefumi 235–I  
 Morita, Shigeyuki 665–I  
 Moscovici, Henri 675–I  
 Murai, Takafumi 901–II  
 Murakami, Haruo 1673–II  
 Neishtadt, Anatoly I. 1271–II  
 Nesterenko, Yuri 447–I  
 Newhouse, Sheldon E. 1285–II  
 Ohsawa, Takeo 913–II  
 Pimsner, Michael V. 979–II  
 Popa, Sorin Teodor 987–II  
 Prasad, Gopal 821–II  
 Preiss, David 923–II  
 Rallis, Stephen 833–II  
 Rees, Mary 1295–II  
 Renegar, James 1595–II  
 Reshetikhin, Nicolai 1373–II  
 Roberts, Paul C. 361–I  
 Rödl, Vojtech 1455–II  
 Roggenkamp, Klaus W. 369–I  
 Saito, Kyoji 931–II  
 Saito, Morihiko 725–I  
 Saper, Leslie 735–I  
 Sarnak, Peter C. 459–I  
 Schapira, Pierre 1187–II  
 Schwarz, Albert 1377–II  
 Segal, Graeme 1387–II  
 Shiota, Tetsuji 473–I  
 Shustin, Eugenii I. 559–I  
 Sibony, Nessim 943–II  
 Sigal, I. M. 1397–II  
 Simpson, Carlos T. 747–I  
 Sinai, Yakov G. 249–I  
 Skandalis, Georges 997–II  
 Slaman, Theodore A. 303–I  
 Steenbrink, Joseph H. M. 569–I  
 Struwe, Michael 1197–II  
 Sunada, Toshikazu 577–I  
 Takasaki, Kanehisa 1205–II  
 Talagrand, Michel 1011–II  
 Tardos, Éva 1467–II  
 Tartar, Luc 1215–II  
 Taylor, Michael E. 1225–II  
 Thomason, Robert W. 381–I  
 Tian, Gang 587–I  
 Tsuchiya, Akihiro 1409–II  
 Turaev, Vladimir G. 689–I  
 Uhlenbeck, Karen 261–I  
 Varchenko, Alexandre 281–I  
 Varopoulos, Nicholas Th. 951–II  
 Vojta, Paul 757–I  
 Volberg, Alexander L. 959–II  
 Wigderson, Avi 1537–II  
 Yor, Marc 1105–II  
 Zelmanov, Efim I. 395–I