# Bioinformatics and Computational Biology in the post-genomic era

## Rita Casadio

*BIOCOMPUTING GROUP*
*University of Bologna, Italy*

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Bologna Computational Biology Network
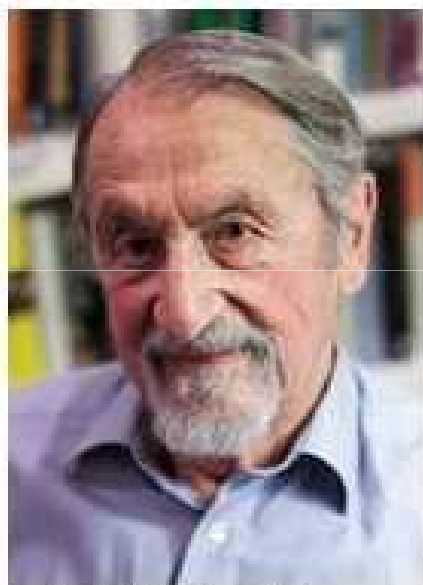
WetLab    CompuLab    SME Small Medium Eterprises

AIRBBC

**Syllabus:**

1) Motivations & definitions
2) The "omic" revolution
3) Next Generation Sequencing Data
4) Data archives  & zooming in on biological complexity
5) Open problems in the omic era
6) Annotation pipelines at the Biocomputing group

# Motivations

## The Nobel Prize in Chemistry 2013

**Martin Karplus**
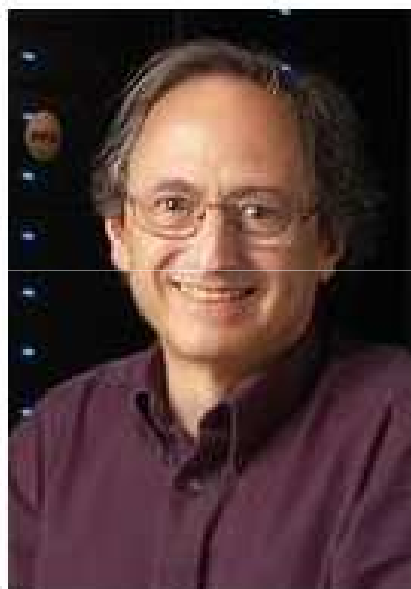© Nobel Media AB

**Michael Levitt**
Photo: Keilana via Wikimedia Commons

**Arieh Warshel**
Photo: Wikimedia Commons

The Nobel Prize in Chemistry 2013 was awarded jointly to Martin Karplus, Michael Levitt and Arieh Warshel *"for the development of multiscale models for complex chemical systems"*.
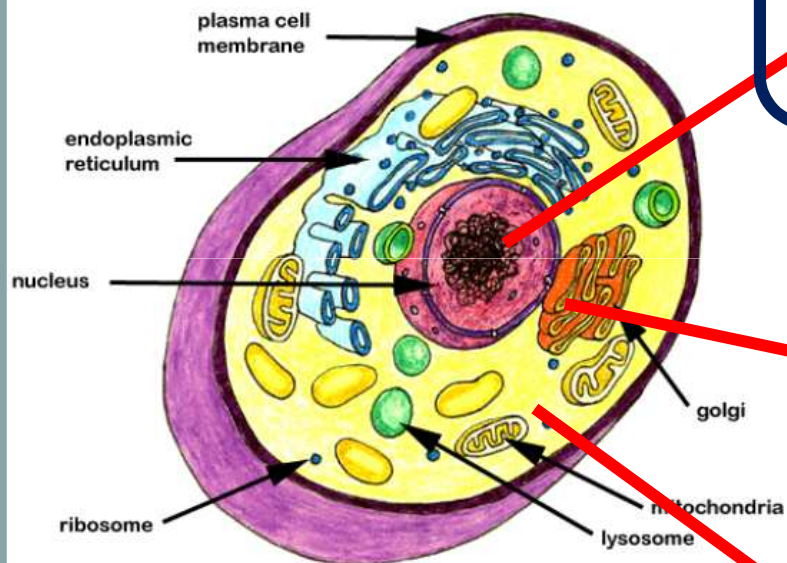
The Nobel Prize in Chemistry 2013 was awarded jointly to Martin Karplus, Michael Levitt and Arieh Warshel *"for the development of multiscale models for complex chemical systems"*.

1. Complex chemical systems
2. Multiscale models
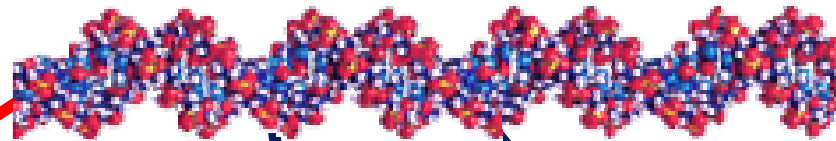3. Development of multiscale models

# The ingredients of biological complexity at the cell level

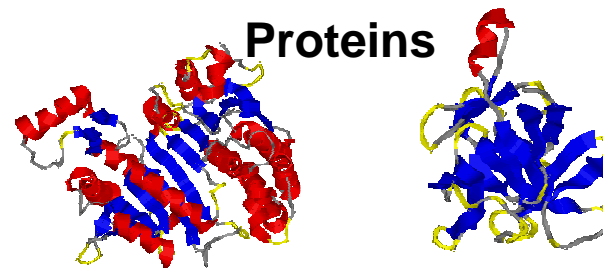*From genes to proteins, their interaction and the interplay with the environment*



**The cell and the environment**

plasma cell membrane

endoplasmic reticulum

nucleus

ribosome

golgi

mitochondria

lysosome

**DNA & Genes**

**Proteins**

**Macromolecular crawding**

*Molecular interactions and functions are affected also by the environment*
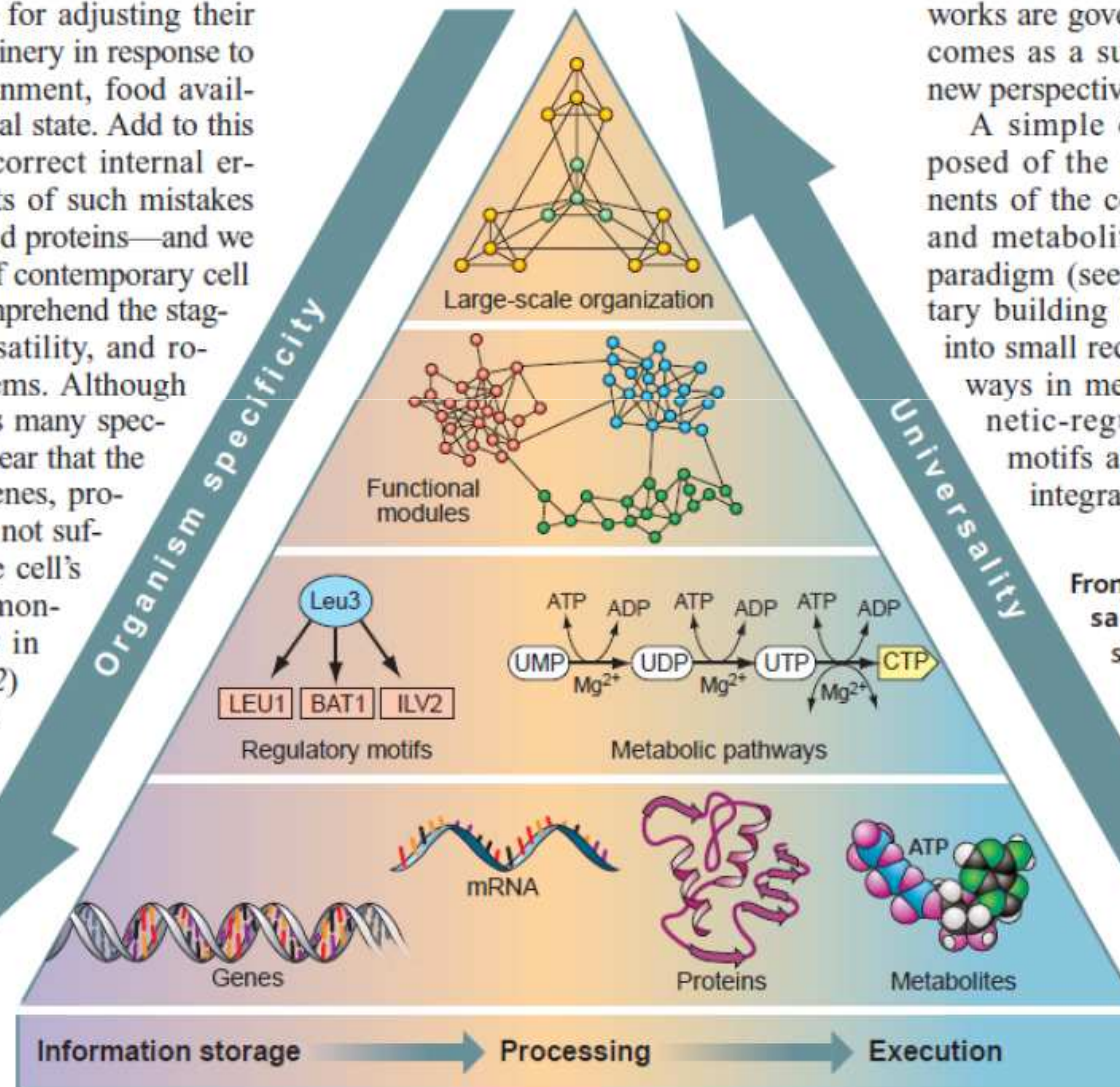
# Life's Complexity Pyramid

Zoltán N. Oltvai and Albert-László Barabási

**SCIENCE VOL 298 ,2002**

Cells and microorganisms have an impressive capacity for adjusting their intracellular machinery in response to changes in their environment, food availability, and developmental state. Add to this an amazing ability to correct internal errors—battling the effects of such mistakes as mutations or misfolded proteins—and we arrive at a major issue of contemporary cell biology: our need to comprehend the staggering complexity, versatility, and robustness of living systems. Although molecular biology offers many spectacular successes, it is clear that the detailed inventory of genes, proteins, and metabolites is not sufficient to understand the cell's complexity (1). As demonstrated by two papers in this issue—Lee et al. (2) on page 799 and Milo et al. (3) on page 824—viewing the cell as a network of genes and proteins offers a viable strategy for addressing the complexity of living systems.

According to the

within large networks (6, 7). evidence for the existence o networks: For example, the p nizes itself into a protein ir work and metabolites are i through an intricate metaboli finding that the structures works are governed by the s comes as a surprise, howev new perspective on cellular o

A simple complexity p posed of the various mole nents of the cell—genes, RI and metabolites—summar paradigm (see the figure). tary building blocks organi into small recurrent pattern ways in metabolism and netic-regulatory netwo motifs and pathways integrated to form fu

**From the particula** sal. The bottom shows the traditi tion of the cell ganization: scriptome, metabol There is tegrat ous the tl



Large-scale organization

Functional modules

Regulatory motifs

Metabolic pathways

Genes — mRNA — Proteins — Metabolites

Organism specificity

Universality

Information storage → Processing → Execution

of cellular c

# An Information Portal to Biological Macromolecular Structures



E.G: RNA Polymerase II Elongation Complex

# The basic information flow: from DNA to proteins

A,T,C,G

cctgttgatggcgacagggactgtatgctgatctatgctgatgcatgcatgctgactactgatgtgggggctat



**From genes...**

>BGAL_SULSO BETA-GALACTOSIDASE Sulfolobus solfataricus
MYSFPNSFRFGWSQAGFQSEMGTPGSEDPNTDWYKWVHDPENMAAGLVSG
DLPENGPGYWGNYKTFHDNAQKMGLKIARLNVEWSRIFPNPLPRPQNFDE
SKQDVTEVEINENELKRLDEYANKDALNHYREIFKDLKSRGLYFILNMYH
WPLPLWLHDPIRVRRGDFTGPSGWLSTRTVYEFARFSAYIAWKFDDLVDE
YSTMNEPNVVGGLGYVGVKSGFPPGYLSFELSRRHMYNIIQAHARAYDGI

A,C,D,E,F,G,H,I,K,L
M,N,P,Q,R,S,T,V,Y,W

**...to Proteins**

# The Data Bases of Biological Sequences and Structures

```
>ENA|M34696|M34696.1 S.solfataricus beta-D-galactosidase (lacS)
gene, complete cds. : Location:1..1000
AAGGAGAAACTTGGCAGTTTATAACTTGACAGTAGGTTGTGGAGTGATGACTGGATCAAT
ACTAGGAGGAGTAGCATATAATTACGTTACACAATTTTATAACCCAATATATTCAATAGA
CCTTATGCTTATCCTATCCTCTATTCTAAGATTCTCGGTATCTCCCCTATTCTTGACCAT
AAAAGATACTCGCTCAAAGCTTAAATAATATTAATCATAAATAAAGTCATGTACTCATTT
CCAAATAGCTTTAGGTTTGGTTGGTCCCAGGCCGGATTTCAATCAGAAATGGGAACACCA
GGGTCAGAAGATCCAAATACTGACTGGTATAAATGGGTTCATGATCCAGAAAACATGGCA
GCGGGATTAGTAAGTGGAGATCTACCAGAAAATGGGCCAGGCTACTGGGGAAACTATAAG
ACATTTCACGATAATGCACAAAAAATGGGATTAAAAATAGCTAGACTAAATGTGGAATGG
TCTAGGATATTTCCTAATCCATTACCAAGGCCACAAAACTTTGATGAATCAAAACAAGAT
GTGACAGAGGTTGAGATAAACGAAAACGAGTTAAAGAGACTTGACGAGTACGCTAATAAA
GACGCATTAAACCATTACAGGGAAATATTCAAGGATCTTAAAAGTAGAGGACTTTACTTT
ATACTAAACATGTATCATTGGCCATTACCTCTATGGTTACACGACCCAATAAGAGTAAGA
AGAGGAGATTTTACTGGACCAAGTGGTTGGCTAAGTACTAGAACAGTTTACGAATTCGCT
AGATTCTCAGCTTATATAGCTTGGAAATTCGATGATCTAGTGGATGAGTACTCAACAATG
AATGAACCTAACGTTGTTGGAGGTTTAGGATACGTTGGTGTTAAGTCCGGTTTTCCCCCA
GGATACCTAAGCTTTGAACTTTCCCGTAGGCATATGTATAACATCATTCAAGCTCACGCA
AGAGCGTATGATGGGATAAAGAGTGTTTCTAAAAAACCAG
```
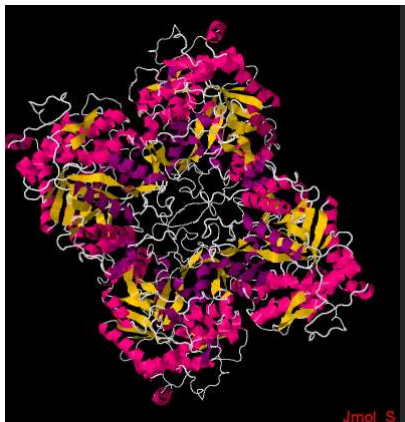
GenBank

```
>BGAL_SULSO BETA-GALACTOSIDASE Sulfolobus solfataricus.
MYSFPNSFRFGWSQAGFQSEMGTPGSEDPNTDWYKWVHDPENMAAGLVSG
DLPENGPGYWGNYKTFHDNAQKMGLKIARLNVEWSRIFPNPLPRPQNFDE
SKQDVTEVEINENELKRLDEYANKDALNHYREIFKDLKSRGLYFILNMYH
WPLPLWLHDPIRVRRGDFTGPSGWLSTRTVYEFARFSAYIAWKFDDLVDE
YSTMNEPNVVGGLGYVGVKSGFPPGYLSFELSRRHMYNIIQAHARAYDGI
KSVSKKPVGIIYANSSFQPLTDKDMEAVEMAENDNRWWFFDAIIRGEITR
GNEKIVRDDLKGRLDWIGVNYYTRTVVKRTEKGYVSLGGYGHGCERNSVS
LAGLPTSDFGWEFFPEGLYDVLTKYWNRYHLMYVTENGIADDADYQRPY
YLVSHVYQVHRAINSGADVRGYLHWSLADNYEWASGFSMRFGLLKVDYNT
KRLYWRPSALVYREIATNGAITDEIEHLNSVPPVKPLRH
```
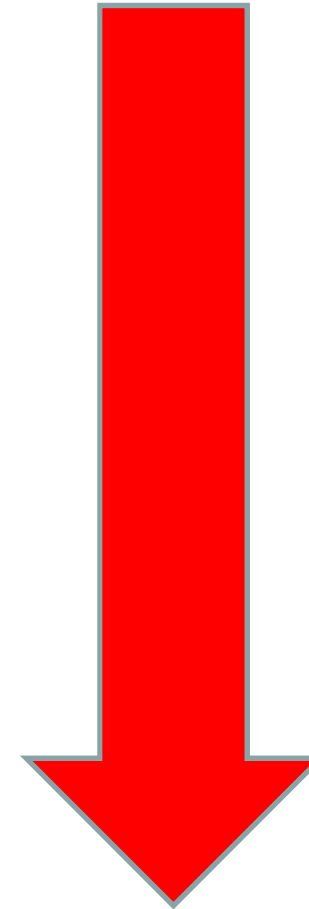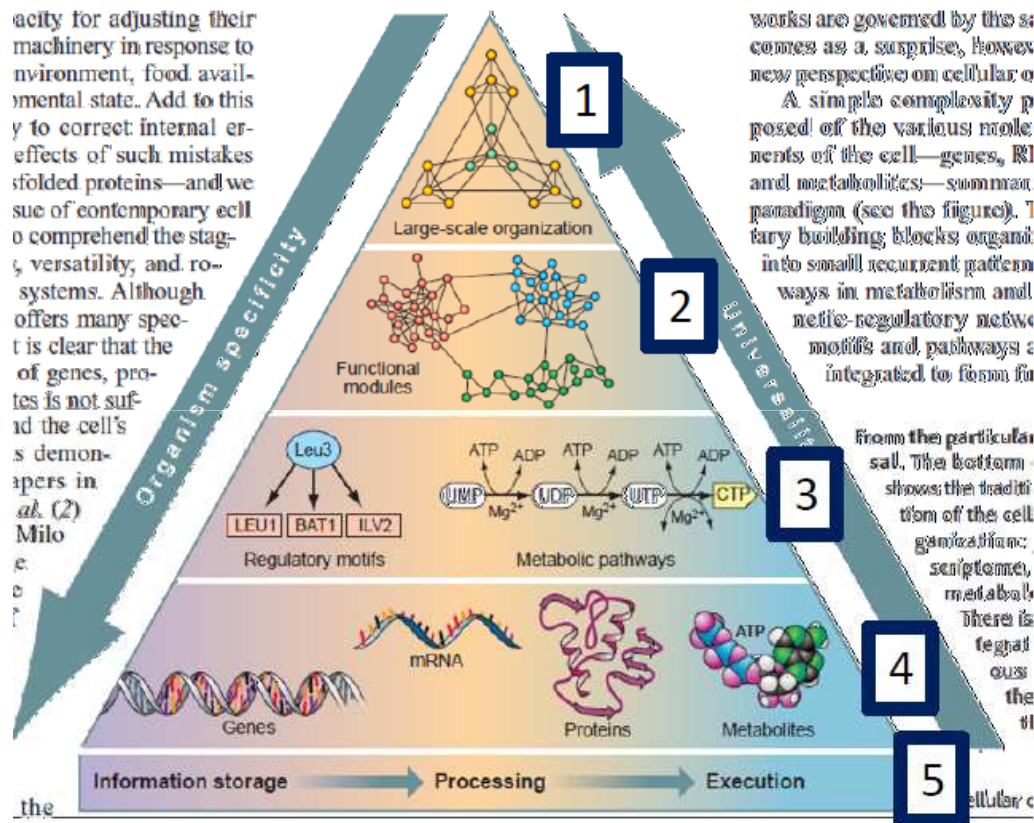
UniProt/SwissProt

PDB

1GOV

INFORMATION − +

# Hierchical levels of cell complexity and our knoweldge



1) **Large-scale organisation**

2) **Functional modules**

3) **Regulatory motifs, metabolic pathways**

4) **Molecules: genes, mRNAs, proteins, metabolites**

5) **Overall: Information storage, Processing, Execution**

# BIOINFORMATICS

## Data Bases

(Biosequences, Structures, Genomes, DNA Chips, Proteomes, Interatomics, Literature)

- Implementation
- Data Mining
- Links

## Computational Biology

Tools for:
- Sequence analysis
- Functional genomics
- Proteomics

## Systems Biology

Models for:
Interatomics, Methabolomics, Evolving complex biosystems (Cell, Organism,..)

*Going back to definitions……who said what and when….*

# NIH WORKING DEFINITION OF BIOINFORMATICS AND COMPUTATIONAL BIOLOGY
## July 17, 2000

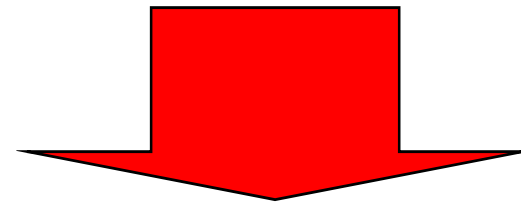The following working definition of bioinformatics and computational biology were developed by the BISTIC Definition Committee and released on July 17, 2000. The committee was chaired by Dr. Michael Huerta of the National Institute of Mental Health and consisted of the following members:

**Bioinformatics Definition Committee**

| **BISTIC Members** | **Expert Members** |
| --- | --- |
| Michael Huerta (Chair) | Gregory Downing |
| Florence Haseltine | Belinda Seto |
| Yuan Liu | |

*BISTIC: Biomedical Information Science and Technology Initiative Consortium*

**Definition**

The NIH Biomedical Information Science and Technology Initiative Consortium agreed on the following definitions of bioinformatics and computational biology recognizing that no definition could completely eliminate overlap with other activities or preclude variations in interpretation by different individuals and organizations.

*Bioinformatics*: Research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data.

*Computational Biology*: The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems.

# HORIZONS

# Life, logic and information

Paul Nurse

**Focusing on information flow will help us to understand better how cells and organisms work.**
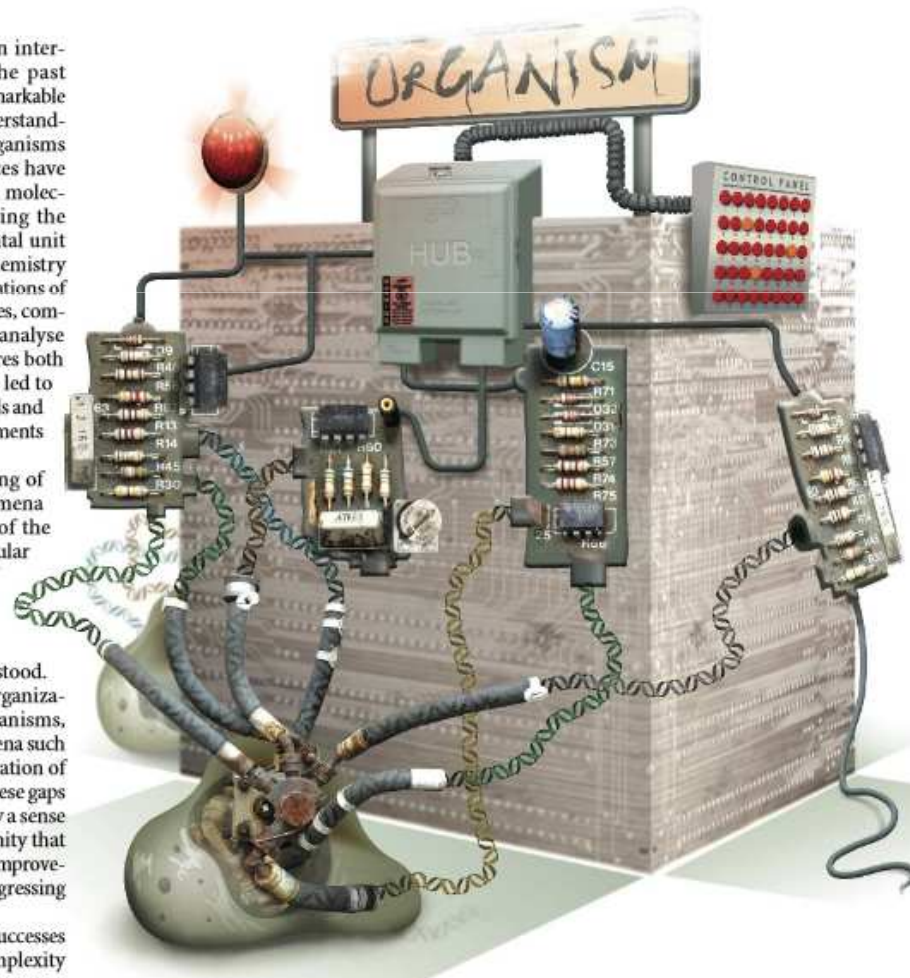


**Some references**

Biology stands at an interesting juncture. The past decades have seen remarkable advances in our understanding of how living organisms work. These advances have been built mostly on molecular biology: applying the ideas that the gene is the fundamental unit of biological information and that chemistry provides effective mechanistic explanations of biological processes. These approaches, combined with an increasing ability to analyse highly complex biomolecular mixtures both qualitatively and quantitatively, have led to our present good understanding of cells and organisms and to significant improvements in our knowledge of human disease.

But comprehensive understanding of many higher-level biological phenomena remains elusive. Even at the level of the cell, phenomena such as general cellular homeostasis and the maintenance of cell integrity, the generation of spatial and temporal order, inter- and intracellular signalling, cell 'memory' and reproduction are not fully understood.

This is also true for the levels of organization seen in tissues, organs and organisms, which feature more complex phenomena such as embryonic development and operation of the immune and nervous systems. These gaps in our knowledge are accompanied by a sense of unease in the biomedical community that understanding of human disease and improvements in disease management are progressing too slowly.

One reason for this is that our past successes have led us to underestimate the complexity

## Perspective

# The Roots of Bioinformatics in Theoretical Biology

**Paulien Hogeweg**<sup></sup>

Theoretical Biology and Bioinformatics Group, Department of Biology, Faculty of Science, Utrecht University, Utrecht, The Netherlands

**Some references**

**Abstract:** From the late 1980s onward, the term "bioinformatics" mostly has been used to refer to computational methods for comparative analysis of genome data. However, the term was originally more widely defined as the study of informatic processes in biotic systems. In this essay, I will trace this early history (from a personal point of view) and I will argue that the original meaning of the term is re-emerging.

## Early History: Bioinformatics, a Work Concept

In the beginning of the 1970s, Ben Hesper and I started to use the term "bioinformatics" for the research we wanted to do, defining it as "the study of

Information" [5] summarized the state of the art in molecular biology before the "sequence age", unraveling for me the essential processes that, at the time in genetics undergraduate texts, were buried in "bead genetics". It seems that recently, after a dormant phase, such information-centric terminology has become more prevalent again (e.g., in terms of identifying a distinct research field [4] and focusing on such processes as sensing the environment [6] and dynamic phosphorylation and methylation codes [7,8]).

We were embedded then within theoretical biology. At the time, after general systems theory [9,10] had come and gone, theoretical biology was in a mild resurgence in acceptance. The series of books entitled "Towards a Theoretical Biology", edited by Waddington [11] (reprints of which are underway), had appeared a few years earlier. In 1972, the main topic at a

enzyme dynamics (e.g., [15,16]), positional information [17], and bi-stability in gene regulation [18] were presented and hotly discussed. Spatial pattern formation was one of the central topics, contrasting Turing systems [19] with gradient-based systems [17]. Francis Crick, who in that period published some papers on gradients in development [20], attended the meeting. Skeptical about the emphasis Turing Patterns were (still) receiving, Crick quoted Turing as saying in reaction to enthusiasm about his work: "Well, the stripes are easy but what about the horse part?" To go "for the horse part", i.e., to go beyond pattern formation to multilevel models of development and morphogenesis, became one of the long-term goals of our nascent work concept "bioinformatics".

Also at about that time, John Maynard Smith gave a lecture in Utrecht and posed a similar challenge with respect to evolu-

# The "omic" revolution

**The analysis of the components of a living organism in its entirety**

Biology becomes a data driven science

GA EVANS, Nature Biotechnology 18:127, 2000

NGS technology allows an unprecedent rate of DNA/RNA sequencing (>4TB per week)

# Next-Generation Sequencing

A large number of platforms using different strategies and chemistries, and with a different throughput are entering the market.

**Results:> 3000 fully sequenced genomes; >1000 human genomes; >10,000 human exoms**

Ion Proton

PacBio

Roche / 454 Genome Sequencer FLX titanium (800 bp, 800 Mb / run)

Illumina / Solexa Genetic Analyzer HiSeq 2000 (150x2 bp, 600 Gb / run)

Applied Biosystems SOLiD 4 System™ (100x2 bp, 400 Gb / run)

# Dealing with genomic data…makes modern biology a BIG science

*Scott D. Kahn*
*Science 331, 728 (2011)*

**The World's Most Powerful Supercomputer Is in China: the Tianhe-2. It's a system developed by China's National University of Defense Technology, and it is capable of running at 33.86 petaflops. (A petaflop is a quadrillion calculations per second.)**

**\*million instructions per second (MIPS)**

### Sequencing Progress vs Compute and Storage
Moore's and Kryder's Laws fall far behind

$10^8$
$10^7$
$10^6$
$10^5$

Legend:
— Microprocessor (MIPS)
— Sequencing (kbases/day)
— Compact HDD storage capacity (MB)

Y-axis: 1, 10, 100, 1000, 10000, 100000, 1000000, 10000000, 100000000

X-axis (Year): 1996, 1997, 1998, 1999, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010

**Fig. 1.** A doubling of sequencing output every 9 months has outpaced and overtaken performance improvements within the disk storage and high-performance computation fields.

# The "omic" era-RESULTS



## Complete Genomes

Prokaryotes: 2975

Eukaryotes:   213

Viruses:   4101

http://www.ncbi.nlm.nih.gov/   Update: May 2014

http://www.1000genomes.org/

Eukaryotic Genome archives

http://www.hugo-international.org
**HUGO**

*Worldwide consortia*

**HUPO**
http://www.hupo.org/

# Proteomics

**the large-scale study of proteins, particularly their expression level, structures and functions**



2-dimensionale separation of proteins

**+ MALDI Mass Spectroscopy**

# LIFE IS CROWDED:
# Macromolecular crowding is under-appreciated



The Crowded Cell: This picture shows an atomically detailed model of the crowded E. coli cytoplasm, including the 50 most abundant macromolecules. RNA is shown as green and yellow. Reprinted from: McGuffee SR, Elcock AH (2010) Diffusion, Crowding & Protein Stability in a Dynamic Molecular Model of the Bacterial Cytoplasm. PLoS Comput Biol 6(3): e1000694.

# What did we learn:

**A shift of paradigm….to describe protein-protein and protein/DNA/RNA interactions**

X        Y        Z        W

**from linear**

**to network models**

A protein is a node characterised by a degree of connections (number of possible interactions or number of other proteins/molecules with which it can interact)

# Interactomics

**In terms of proteomics, interactomics refers to protein-protein interaction networks (or protein/DNA/RNA interactions)**



Vidal, M et al. Nature 2005. 437: 1173–1178,

# Search for human genetic variability.....

## Where to check variations for disease association



OMIM Statistics for January 7, 2013

Number of Entries

| | Autosomal | X-Linked | Y-Linked | Mitochondrial | Total |
|---|---|---|---|---|---|
| * Gene with known sequence | 13370 | 651 | 48 | 35 | 14104 |
| + Gene with known sequence and phenotype | 124 | 4 | 0 | 2 | 130 |
| # Phenotype description, molecular basis known | 3371 | 271 | 4 | 28 | 3674 |
| % Mendelian phenotype or locus, molecular basis unknown | 1627 | 133 | 5 | 0 | 1765 |
| Other, mainly phenotypes with suspected mendelian basis | 1765 | 125 | 2 | 0 | 1892 |
| Total | 20257 | 1184 | 59 | 65 | 21565 |

http://www.ncbi.nlm.nih.gov/Omim/mimstats.html

## Where to find disease associated variations in proteins



http://www.uniprot.org/docs/humsavar

## Where to search/check for neutral variations



http://www.ncbi.nlm.nih.gov/projects/SNP/

## The Human Variome

http://www.ornl.gov/hgmis.



# SNPs: Single Nucleotide Polymorphisms

**(about 20,876 genes and 181,744 transcripts in the human genome)**

>protein kinase

acctgttgatggcgacagggactgtatgctgatct
atgctgatgcatgcatgctgactactgatgtggggg
gctattgacttgatgtctatc....

*Genes in DNA...*

*...with different effects depending on variability*

**Over 50 millions of single mutations are known**

*...code for proteins...*

# Overall: from Genotype to Phenotype

*...proteins correspond to functions...*

*....in methabolic pathways*

From 5000 to 10000 proteins per tissue

*Proteins interact*

*...when they are expressed*

**Beyond the post-genomic era**

# Encyclopedia of DNA Elements

**About ENCODE Data**

The Encyclopedia of DNA Elements (ENCODE) Consortium is an international collaboration of researc
Research Institute (NHGRI). The goal of ENCODE is to build a comprehensive parts list of functional ele
act at the protein and RNA levels, and regulatory elements that control cells and circumstances in which

ENCODE data are now available for the en
*available for immediate use via :*

- Search for displayable tracks and downloadable files

**Human**
Integrative Analysis
Experiment Matrix
Experiment List
Search
Downloads
Genome Browser (hg19)
Session Gallery
Cell Types

**Mouse**
Experiment Matrix
Experiment List
Search
Downloads

To sea
feature

All EN
Data F
UCSC

**News**

**Beyond the post-genomic era**



日本語

Home    Data    Views    Protocols    Software    Papers    FAQ

## FANTOM5 papers have been published!

Using Cap Analysis of Gene Expression (CAGE) we have mapped the sets of transcripts, transcription factors, promoters and enhancers active in the majority of mammalian primary cell types. We have also complemented this with profiles from cancer cell lines, and tissues. The results are described in two landmark papers in Nature describing the promoterome and enhancerome of mammalian cells. An additional 16 publications cover areas as diverse as primary cells, gene families, genome wide observations on promoter features and new bioinformatics tools.

## FANTOM

FANTOM is an international research consortium established by Dr. Hayashizaki and his colleagues in 2000 to assign functional annotations to the full-length cDNAs that were collected during the Mouse Encyclopedia Project at RIKEN. FANTOM has since developed and expanded over time to encompass the fields of transcriptome analysis. The object of the project is moving steadily up the layers in the system of life, progressing thus from an understanding of the 'elements' - the transcripts - to an understanding of the 'system' - the transcriptional regulatory network, in other words the 'system' of an individual life form.

## LATEST NEWS

**Apr 30, 2014** Two new FANTOM5 related publications
Read More »

# Functional Genomics/Epigenomics

## *From genes to functions and backward*

**DATA -INTEGRATION**

The "omic" era

Complexity

- Genomics
- Transcriptomics
- Proteomics
- Metabolomics
- Regulomics

Systems Biology

**Summing up…..**

**Open problems in the post-genomic era after DNA/RNA sequencing**

1) Genome assembly & Genome annotation (e.g. exon/intron boundaries)
2) Chromatin dynamics
3) Finding alternative splicing variants
4) Protein structural and functional annotation
5) Annotation of SNP variants & Correlation among SNPs and diseases
6) Simulation of cell complexity

# Genomic data and the problem of protein validation

## Data production→Data analysis

DNA sequencing →gene recognition → protein translation

**Experiments to validate protein structure and function produce data in a time >> than that required to deposit putative protein sequences into data bases**

# A "BIG" problem of the "omic era" after genome sequencing:



code for

that are endowed with

```
        10        20        30        40        50        60        70
        |         |         |         |         |         |         |
TEKLWVTVYY GVPVWKEATT TLFCASDAKA YDTEVHNVWA THACVPTDPN PQEVVLVNVT ENFNMWKNDM
        80        90       100       110       120       130       140
        |         |         |         |         |         |         |
VEQMHEDIIS LWDQSLKPCV KLTPLCVSLK CTDLKNDTNT NSSSGRMIME KGEIKNCSFN ISTSIRGKVQ
       150       160       170       180       190       200       210
        |         |         |         |         |         |         |
KEYAFFYKLD IIPIDNDTTS YKLTSCNTSV ITQACPKVSF EPIPIHYCAP AGFAILKCNN KTFNGTGPCT
       220       230       240       250       260       270       280
        |         |         |         |         |         |         |
NVSTVQCTHG IRPVVSTQLL LNGSLAEEEV VIRSVNFTDN AKTIIVQLNT SVEINCTRPN NNTRKRIRIQ
       290       300       310       320       330       340       350
        |         |         |         |         |         |         |
RGPGRAFVTI GKIGNMRQAH CNISRAKWNN TLKQIASKLR EQFGNNKTII FKQSSGGDPE IVTHSFNCGG
       360       370       380       390       400       410       420
        |         |         |         |         |         |         |
EFFYCNSTQL FNSTWFNSTW STEGSNNTEG SDTITLPCRI KQIINMWQKV GKAMYAPPIS GQIRCSSNIT
       430       440       450       460       470       480
        |         |         |         |         |         |
GLLLTRDGGN SNNESEIFRP GGGDMRDNWR SELYKYKVVK IEPLGVAPTK AKRRVVQREK R
```

**Protein sequences (~56millions)**



**Genomes (~3000)**



**Protein structures and functions**

## Protein sequence Annotation:
## to endow with structural and functional features protein sequences after gene translation

# Protein function

Open menus

**The Ontologies**
- Cellular component
- Biological process
- Molecular function

## GO function vocabulary:
http://www.geneontology.org/

## Ontology Structure

The Gene Ontology is a **controlled vocabulary**, a set of standard terms—words and phrases—used for indexing and retrieving information. In addition to defining terms, GO also defines the **relationships** between the terms, making it a **structured** vocabulary.

## GO as a Graph

The structure of GO can be described in terms of a graph, where each GO term is a node, and the relationships between the terms are arcs between the nodes. The relationships used in GO are **directed**—for example, a mitochondrion *is an* organelle, but an organelle is not a mitochondrion—and the graph is **acyclic**, meaning that cycles are not allowed in the graph. The ontologies resemble a hierarchy, as child terms are more specialized and parent terms are less specialized, but unlike a hierarchy, a term may have more than one parent term. For example, the biological process term hexose biosynthetic process has two parents, hexose metabolic process and monosaccharide biosynthetic process. This is because biosynthetic process is a type of metabolic process and a hexose is a type of monosaccharide.

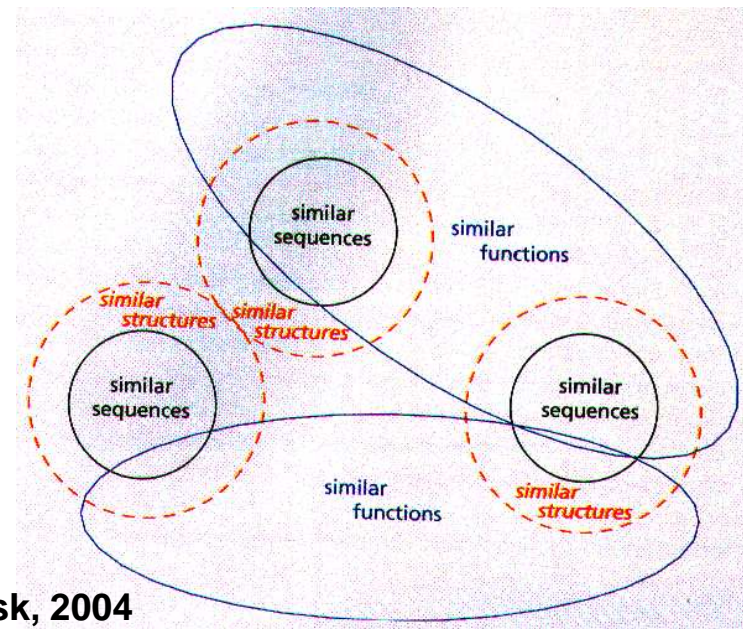# Protein annotation by sequence similarity

```
Homology between CT46 and MGC26710 hypothetical protein

Identities = 136/249 (54%), with conservative changes = 180/249 (72%)

CT46       1    MATAQLQR-----TPMSALVFPNKISTEHQSLVLVKRLLAVSVSCITYLRGIFPECAYGTRYLDDLCVKILREDK
                MATAQL       VFP++I+ EH+SL +VK+L A S+SCITYLRG+FPE +YG R+LDDL +KILREDK
MGC26710   1    MATAQLSHCITIHKASKETVFPSQITNEHESLKMVKKLFATSISCITYLRGLFPESSYGERHLDDLSLKILREDK

CT46       71   NCPGSTQLVKWMLGCYDALQKKYLRMVVLAVYTNPEDPQTISECYQFKFKYTNNGPLMDF--ISKNQSNESSMLS
                 CPGS  +++W+ GC+DAL+K+YLRM VL +YT+P   + ++E YQFKFKYT  G  MDF   S + S ES   +
MGC26710   76   KCPGSLHIIRWIQGCFDALEKRYLRMAVLTLYTDPMGSEKVTEMYQFKFKYTKEGATMDFDSHSSSTSFESGTNN

CT46       144  TDTKKASILLIRKIYILMQNLGPLPNDVCLTMKLFYYDEVTPPDYQPPGFKDG-DCEGVIFEGEPMYLNVGEVST
                 D KKAS+LLIRK+YILMQ+L PLPN+V LTMKL YY+ VTP DYQP GFK+G +   ++F+ EP+ + VG VST
MGC26710   151  EDIKKASVLLIRKLYILMQDLEPLPNNVVLTMKLHYYNAVTPHDYQPLGFKEGVNSHFLLFDKEPINVQVGFVST

CT46       218  PFHIFKVKVTTERERMENIDSTIL  241
                 FH  KVKV TE  ++ ++++ +
MGC26710   226  GFHSMKVKVMTEATKVIDLENNLF  249
```

**If sequences share more than 30/40% sequence identity they can share similar structure and function**

*But the problem is much more complex*



**Lesk, 2004**

# Transfer of annotation *in silico* by homology search

```
ADH1_SULSO   ----------MRAVRLVEIGKP--LSLQEIGVPKPKGPQVLIKVEAAGVCHSDVHMRQGRFGNLRIVE
ADH_CLOBE    ----------MKGFAMLGINKLG---WIEKERPVAGSYDAIVRPLAVSPCTSDIHTVFEGA-------
ADH_THEBR    ----------MKGFAMLSIGKVG---WIEKEKPAPGPFDAIVRPLAVAPCTSDIHTVFEGA-------
ADH1_SOLTU   MSTTVGQVIRCKAAVAWEAGKP--LVMEEVDVAPPQKMEVRLKILYTSLCHTDVYFWEAKG-------
ADH2_LYCES   MSTTVGQVIRCKAAVAWEAGKP--LVMEEVDVAPPQKMEVRLKILYTSLCHTDVYFWEAKG-------
ADH1_ASPFL   ----MSIPEMQWAQVAEQKGGP--LIYKQIPVPKPGPDEILVKVRYSGVCHTDLHALKGDW-------
```

**Sequence comparison is performed with alignment programs**

Sequence identity ≥ 30 %  ⟹  3D ?; Similar function ??

# Methods for similarity searches:

**BLAST, Psi-BLAST (http://www.ncbi.nlm.nih.gov/BLAST/)**

**Altschul et al., (1990)  J Mol Biol 215:403-410**

**Altschul et al., (1998) Nucleic Acids Res. 25:3389-3402**
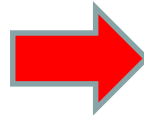
**Pfam (http://pfam.wustl.edu/hmmsearch.shtml)**

**Bateman et al., (2000) Nucleic Acids Research 28:263-266**

**The little we know (SwissProt)……is expanded to annotate all the protein sequences (TrEMBL)**

*Release 2014_03 (19-Mar) of UniProtKB/Swiss-Prot contains 542,782 sequence entries:*

| Protein existence (PE): | entries | % |
|---|---|---|
| 1: Evidence at protein level | 82,087 | 15.1 |
| 2: Evidence at transcript level | 62,227 | 11.5 |
| 3: Inferred from homology | 380,832 | 70.2 |
| 4: Predicted | 15,705 | 2.9 |

*Release 2014_03 (19-Mar) of UniProtKB/TrEMBL contains 54,247,468 sequence entries:*

| Protein existence (PE): | entries | % |
|---|---|---|
| 1: Evidence at protein level | 22,013 | 0.04 |
| 2: Evidence at transcript level | 931,313 | 1.72 |
| 3: Inferred from homology | 13,573,938 | 25.02 |
| 4: Predicted | 39,720,204 | 73.22 |

**Automatic annotation at UNIPROTKB**
*http://www.uniprot.org/program/automatic_annotation*

UniProt has developed two prediction systems, **UniRule** and **the Statistical Automatic Annotation System (SAAS)** to automatically annotate UniProtKB/TrEMBL in an efficient and scalable manner with a high degree of accuracy:

- **Based on rules**
- **Rules are created, tested and validated against published experimental data in UniProtKB/Swiss-Prot**
- **Rules are linked to InterPro member database signatures**
- **Rules have annotations and conditions**
- **Rules are reapplied to UniProtKB/TrEMBL every four-weekly release with both automatic and manual QA procedures ensuring they are still valid**