



GERHARD WEIKUM  
GOTTFRIED VOSSEN

# TRANSACTIONAL INFORMATION SYSTEMS

THEORY, ALGORITHMS, AND THE PRACTICE  
OF CONCURRENCY CONTROL AND RECOVERY

# **Transactional Information Systems**

Theory, Algorithms, and the Practice of  
Concurrency Control and Recovery

## **The Morgan Kaufmann Series in Data Management Systems**

*Series Editor:* Jim Gray, Microsoft Research

*Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery*

Gerhard Weikum and Gottfried Vossen

*Information Visualization in Data Mining and Knowledge Discovery*

Edited by Usama Fayyad, Georges G. Grinstein, and Andreas Wierse

*Information Modeling and Relational Databases: From Conceptual Analysis to Logical Design*

Terry Halpin

*Spatial Databases: With Application to GIS*

Philippe Rigaux, Michel Scholl, and Agnes Voisard

*SQL: 1999—Understanding Relational Language Components*

Jim Melton and Alan R. Simon

*Component Database Systems*

Edited by Klaus R. Dittrich and Andreas Geppert

*Managing Reference Data in Enterprise Databases: Binding Corporate Data to the Wider World*

Malcolm Chisholm

*Data Mining: Concepts and Techniques*

Jiawei Han and Micheline Kamber

*Understanding SQL and Java Together: A Guide to SQLJ, JDBC, and Related Technologies*

Jim Melton and Andrew Eisenberg

*Database: Principles, Programming, and Performance*, Second Edition

Patrick and Elizabeth O'Neil

*The Object Data Standard: ODMG 3.0*

Edited by R. G. G. Cattell and Douglas K. Barry

*Data on the Web: From Relations to Semistructured Data and XML*

Serge Abiteboul, Peter Buneman, and Dan Suciu

*Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*

Ian H. Witten and Eibe Frank

*Joe Celko's SQL for Smarties: Advanced SQL Programming*, Second Edition

Joe Celko

*Joe Celko's Data and Databases: Concepts in Practice*

Joe Celko

*Developing Time-Oriented Database Applications in SQL*

Richard T. Snodgrass

*Web Farming for the Data Warehouse*

Richard D. Hackathorn

*Database Modeling & Design*, Third Edition

Toby J. Teorey

*Management of Heterogeneous and Autonomous Database Systems*

Edited by Ahmed Elmagarmid, Marek Rusinkiewicz, and Amit Sheth

*Object-Relational DBMSs: Tracking the Next Great Wave*, Second Edition

Michael Stonebraker and Paul Brown, with Dorothy Moore

*A Complete Guide to DB2 Universal Database*

Don Chamberlin

*Universal Database Management: A Guide to Object/Relational Technology*

Cynthia Maro Saracco

*Readings in Database Systems*, Third Edition

Edited by Michael Stonebraker and Joseph M. Hellerstein

*Understanding SQL's Stored Procedures: A Complete Guide to SQL/PSM*

Jim Melton

*Principles of Multimedia Database Systems*

V. S. Subrahmanian

*Principles of Database Query Processing for Advanced Applications*

Clement T. Yu and Weiyi Meng

*Advanced Database Systems*

Carlo Zaniolo, Stefano Ceri, Christos Faloutsos, Richard T. Snodgrass,

V. S. Subrahmanian, and Roberto Zicari

*Principles of Transaction Processing: For the Systems Professional*

Philip A. Bernstein and Eric Newcomer

*Using the New DB2: IBM's Object-Relational Database System*

Don Chamberlin

*Distributed Algorithms*

Nancy A. Lynch

*Active Database Systems: Triggers and Rules for Advanced Database Processing*

Edited by Jennifer Widom and Stefano Ceri

*Migrating Legacy Systems: Gateways, Interfaces, & the Incremental Approach*

Michael L. Brodie and Michael Stonebraker

*Atomic Transactions*

Nancy Lynch, Michael Merritt, William Weihl, and Alan Fekete

*Query Processing for Advanced Database Systems*

Edited by Johann Christoph Freytag, David Maier, and Gottfried Vossen

*Transaction Processing: Concepts and Techniques*

Jim Gray and Andreas Reuter

*Building an Object-Oriented Database System: The Story of O<sub>2</sub>*

Edited by François Bancilhon, Claude Delobel, and Paris Kanellakis

*Database Transaction Models for Advanced Applications*

Edited by Ahmed K. Elmagarmid

*A Guide to Developing Client/Server SQL Applications*

Setrag Khoshafian, Arvola Chan, Anna Wong, and Harry K. T. Wong

*The Benchmark Handbook for Database and Transaction Processing Systems*,

Second Edition

Edited by Jim Gray

*Camelot and Avalon: A Distributed Transaction Facility*

Edited by Jeffrey L. Eppinger, Lily B. Mummert, and Alfred Z. Spector

*Readings in Object-Oriented Database Systems*

Edited by Stanley B. Zdonik and David Maier

This Page Intentionally Left Blank

# Transactional Information Systems

Theory, Algorithms, and the Practice of  
Concurrency Control and Recovery

Gerhard Weikum

*University of the Saarland, Germany*

Gottfried Vossen

*University of Münster, Germany*



**MORGAN KAUFMANN PUBLISHERS**

AN IMPRINT OF ACADEMIC PRESS

A Harcourt Science and Technology Company

SAN FRANCISCO SAN DIEGO NEW YORK BOSTON  
LONDON SYDNEY TOKYO

<i>Executive Editor</i>	Diane D. Cerra
<i>Publishing Services Manager</i>	Scott Norton
<i>Assistant Publishing Services Manager</i>	Edward Wade
<i>Assistant Editor</i>	Belinda Breyer
<i>Cover Design</i>	Frances Baca Design
<i>Cover Image</i>	© Ralph A. Clevenger/CORBIS
<i>Text Design</i>	Rebecca Evans & Associates
<i>Composition</i>	TechBooks
<i>Technical Illustration</i>	Dartmouth Publishing, Inc.
<i>Copyeditor</i>	Judith Brown
<i>Proofreader</i>	Jennifer McClain
<i>Indexer</i>	Steve Rath
<i>Printer</i>	Courier Corporation

Designations used by companies to distinguish their products are often claimed as trademarks or registered trademarks. In all instances in which Morgan Kaufmann Publishers is aware of a claim, the product names appear in initial capital or all capital letters. Readers, however, should contact the appropriate companies for more complete information regarding trademarks and registration.

Morgan Kaufmann Publishers  
340 Pine Street, Sixth Floor, San Francisco, CA 94104-3205, USA  
<http://www.mkp.com>

ACADEMIC PRESS  
A Harcourt Science and Technology Company  
525 B Street, Suite 1900, San Diego, CA 92101-4495, USA  
<http://www.academicpress.com>

Academic Press  
Harcourt Place, 32 Jamestown Road, London, NW1 7BY, United Kingdom  
<http://www.academicpress.com>

© 2002 by Academic Press  
All rights reserved  
Printed in the United States of America

06 05 04 03 02 5 4 3 2 1

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopying, or otherwise—without the prior written permission of the publisher.

**Library of Congress Cataloging-in-Publication Data**

This book is printed on acid-free paper.

To science  
and for Liz, Maria, and Julia, my mother  
Anna, and my late father Oskar.

—*Gerhard Weikum*

For Martina, for the first 24 years and for staying by my side  
through the completion of many books, for Laura and  
Kathrin, my true challenges, for my mother Marianne,  
our sponsor, and for my father, Hans, whose life  
transaction was serialized by heaven's scheduler before  
the transaction of finishing this book, as he had to  
leave this world about nine hours before the central  
shadow of the 1999 solar eclipse hit Germany  
(and I firmly believe it was a Commit).

—*Gottfried Vossen*

Science! True daughter of Old Time thou art!  
Who alterest all things with thy peering eyes.

—*Edgar Allan Poe*



This Page Intentionally Left Blank

# Foreword

Jim Gray, Microsoft, Inc.

This book is a major advance for transaction processing. It synthesizes and organizes the last three decades of research into a rigorous and consistent presentation. It unifies concurrency control and recovery for both the page and object models. As the copious references show, this unification has been the labor of many researchers in addition to Weikum and Vossen; but this book organizes that huge research corpus into a consistent whole, with a step-by-step development of the ideas.

The “classic” books on transaction processing have largely either focused on the practical aspects or taken a rigorous approach presenting theorems and proofs. Most have limited themselves to “flat” transactions because the theory of multilevel transactions was so immature. This is the first book to give an in-depth presentation of both the theory and the practical aspects of the field, and the first to present our new understanding of multilevel (object model) transaction processing.

In reading the book, I was impressed at how much our field has advanced, and how once-complex ideas are now simply explained once the terminology is rationalized, and once the proper perspective is set. You will find it possible to read this book at the superficial level: just following the text, the examples, the definitions, and the theorems. You will also be able to dive as deep as you like into the detailed presentation of the results – both the proofs and the programs. In reviewing the book, I took both perspectives: trying to use it as a reference by diving into the middle of some chapter and seeing how quickly I could find the answer to my question. I also took the linear approach of reading the book. In both cases, the book was very informative and very accessible.

This book is likely to become the standard reference in our field for many years to come.

This Page Intentionally Left Blank

# Contents

<b>Foreword</b>	ix
Jim Gray, Microsoft, Inc.	
<b>Preface</b>	xxi
 <b>PART ONE</b>	
BACKGROUND AND MOTIVATION	
<b>Chapter 1</b> What Is It All About?	3
1.1 Goal and Overview	3
1.2 Application Examples	4
1.2.1 Online Transaction Processing: Debit/Credit Example	5
1.2.2 Electronic Commerce Example	9
1.2.3 Workflow Management: Travel Planning Example	12
1.3 System Paradigms	16
1.3.1 Three-Tier and Two-Tier Architectures	16
1.3.2 Federations of Servers	20
1.4 Virtues of the Transaction Concept	22
1.4.1 Transaction Properties and the Transaction Programming Interface	22
1.4.2 Requirements on Transactional Servers	26
1.5 Concepts and Architecture of Database Servers	27
1.5.1 Architectural Layers of Database Systems	27
1.5.2 How Data Is Stored	30
1.5.3 How Data Is Accessed	32
1.5.4 How Queries and Updates Are Executed	35
1.6 Lessons Learned	37
Exercises	38
Bibliographic Notes	38

<b>Chapter 2</b>	<b>Computational Models</b>	<b>41</b>
2.1	Goal and Overview	41
2.2	Ingredients	42
2.3	The Page Model	43
2.4	The Object Model	47
2.5	Road Map of the Book	53
2.6	Lessons Learned	56
	Exercises	56
	Bibliographic Notes	57

## **PART TWO**

### **CONCURRENCY CONTROL**

<b>Chapter 3</b>	<b>Concurrency Control: Notions of Correctness for the Page Model</b>	<b>61</b>
3.1	Goal and Overview	61
3.2	Canonical Concurrency Problems	62
3.3	Syntax of Histories and Schedules	65
3.4	Correctness of Histories and Schedules	71
3.5	Herbrand Semantics of Schedules	73
3.6	Final State Serializability	76
3.7	View Serializability	82
3.7.1	View Equivalence and the Resulting Correctness Criterion	83
3.7.2	On the Complexity of Testing View Serializability	86
3.8	Conflict Serializability	92
3.8.1	Conflict Relations	93
3.8.2	Class CSR	94
3.8.3	Conflicts and Commutativity	99
3.8.4	Restrictions of Conflict Serializability	101
3.9	Commit Serializability	105
3.10	An Alternative Correctness Criterion: Interleaving Specifications	108
3.11	Lessons Learned	119
	Exercises	120
	Bibliographic Notes	121

<b>Chapter 4</b>	<b>Concurrency Control Algorithms</b>	<b>125</b>
4.1	Goal and Overview	125
4.2	General Scheduler Design	126
4.3	Locking Schedulers	130
4.3.1	Introduction	130
4.3.2	The Two-Phase Locking Protocol	133
4.3.3	Deadlock Handling	138
4.3.4	Variants of 2PL	142
4.3.5	Ordered Sharing of Locks	144
4.3.6	Altruistic Locking	150
4.3.7	Non-Two-Phase Locking Protocols	155
4.3.8	On the Geometry of Locking	162
4.4	Nonlocking Schedulers	166
4.4.1	Timestamp Ordering	166
4.4.2	Serialization Graph Testing	168
4.4.3	Optimistic Protocols	170
4.5	Hybrid Protocols	175
4.6	Lessons Learned	179
	Exercises	180
	Bibliographic Notes	182
<b>Chapter 5</b>	<b>Multiversion Concurrency Control</b>	<b>185</b>
5.1	Goal and Overview	185
5.2	Multiversion Schedules	186
5.3	Multiversion Serializability	189
5.3.1	Multiversion View Serializability	189
5.3.2	Testing Membership in MVSR	193
5.3.3	Multiversion Conflict Serializability	197
5.4	Limiting the Number of Versions	201
5.5	Multiversion Concurrency Control Protocols	203
5.5.1	The MVTO Protocol	203
5.5.2	The MV2PL Protocol	205
5.5.3	The MVSGT Protocol	209
5.5.4	A Multiversion Protocol for Read-Only Transactions	211
5.6	Lessons Learned	213
	Exercises	214
	Bibliographic Notes	215

<b>Chapter 6</b>	<b>Concurrency Control on Objects:</b>	
	Notions of Correctness	217
6.1	Goal and Overview	217
6.2	Histories and Schedules	218
6.3	Conflict Serializability for Flat Object Transactions	223
6.4	Tree Reducibility	228
6.5	Sufficient Conditions for Tree Reducibility	233
6.6	Exploiting State Based Commutativity	240
6.7	Lessons Learned	246
	Exercises	247
	Bibliographical Notes	250
<b>Chapter 7</b>	<b>Concurrency Control Algorithms on Objects</b>	251
7.1	Goal and Overview	251
7.2	Locking for Flat Object Transactions	251
7.3	Layered Locking	252
7.4	Locking on General Transaction Forests	259
7.5	Hybrid Algorithms	265
7.6	Locking for Return Value Commutativity and Escrow Locking	267
7.7	Lessons Learned	271
	Exercises	272
	Bibliographic Notes	274
<b>Chapter 8</b>	<b>Concurrency Control on Relational Databases</b>	277
8.1	Goal and Overview	277
8.2	Predicate-Oriented Concurrency Control	278
8.3	Relational Update Transactions	285
8.3.1	Syntax and Semantics	285
8.3.2	Commutativity and Simplification Rules	287
8.3.3	Histories and Final State Serializability	288
8.3.4	Conflict Serializability	291
8.3.5	Extended Conflict Serializability	293
8.3.6	Serializability in the Presence of Functional Dependencies	295
8.3.7	Summary	298
8.4	Exploiting Transaction Program Knowledge	299
8.4.1	Motivating Example	299
8.4.2	Transaction Chopping	301
8.4.3	Applicability of Chopping	306

8.5	Lessons Learned	308
	Exercises	308
	Bibliographic Notes	311
<b>Chapter 9</b>	<b>Concurrency Control on Search Structures</b>	<b>313</b>
9.1	Goal and Overview	313
9.2	Implementation of Search Structures by B <sup>+</sup> Trees	315
9.3	Key Range Locking at the Access Layer	320
9.4	Techniques for the Page Layer	327
9.4.1	Lock Coupling	328
9.4.2	Link Technique	337
9.4.3	Giveup Technique	339
9.5	Further Optimizations	340
9.5.1	Deadlock-Free Page Latching	340
9.5.2	Enhanced Key Range Concurrency	341
9.5.3	Reduced Locking Overhead	343
9.5.4	Exploiting Transient Versioning	344
9.6	Lessons Learned	344
	Exercises	345
	Bibliographic Notes	347
<b>Chapter 10</b>	<b>Implementation and Pragmatic Issues</b>	<b>349</b>
10.1	Goal and Overview	349
10.2	Data Structures of a Lock Manager	349
10.3	Multiple Granularity Locking and Dynamic Escalation	352
10.4	Transient Versioning	354
10.5	Nested Transactions for Intra-transaction Parallelism	357
10.6	Tuning Options	359
10.6.1	Manual Locking	359
10.6.2	SQL Isolation Levels	360
10.6.3	Short Transactions	364
10.6.4	Limiting the Level of Multiprogramming	367
10.7	Overload Control	369
10.7.1	Feedback-Driven Method	369
10.7.2	Wait-Depth Limitation	373
10.8	Lessons Learned	374
	Exercises	375
	Bibliographic Notes	375



**PART THREE****RECOVERY**

<b>Chapter 11</b>	<b>Transaction Recovery</b>	<b>379</b>
11.1	Goal and Overview	379
11.2	Expanded Schedules with Explicit Undo Operations	381
11.2.1	Intuition and Overview of Concepts	381
11.2.2	The Formal Model	382
11.3	Correctness Criteria for the Page Model	385
11.3.1	Expanded Conflict Serializability	385
11.3.2	Reducibility and Prefix Reducibility	387
11.4	Sufficient Syntactic Conditions	390
11.4.1	Recoverability	391
11.4.2	Avoiding Cascading Aborts	391
11.4.3	Strictness	393
11.4.4	Rigorousness	393
11.4.5	Log Recoverability	398
11.5	Page Model Protocols for Schedules with Transaction Aborts	402
11.5.1	Extending Two-Phase Locking for Strictness and Rigorousness	402
11.5.2	Extending Serialization Graph Testing for Log Recoverability	403
11.5.3	Extending Other Protocols for Log Recoverability	406
11.6	Correctness Criteria for the Object Model	407
11.6.1	Aborts in Flat Object Schedules	407
11.6.2	Complete and Partial Aborts in General Object Model Schedules	416
11.7	Object Model Protocols for Schedules with Transaction Aborts	419
11.8	Lessons Learned	420
	Exercises	421
	Bibliographic Notes	423
<b>Chapter 12</b>	<b>Crash Recovery: Notion of Correctness</b>	<b>427</b>
12.1	Goal and Overview	427
12.2	System Architecture and Interfaces	430
12.3	System Model	434
12.4	Correctness Criterion	437
12.5	Road Map of Algorithms	439

12.6 Lessons Learned	444
Exercises	444
Bibliographic Notes	445
<b>Chapter 13</b> Page Model Crash Recovery Algorithms	447
13.1 Goal and Overview	447
13.2 Basic Data Structures	449
13.3 Redo-Winners Paradigm	453
13.3.1 Actions during Normal Operation	454
13.3.2 Simple Three-Pass Algorithm	458
13.3.3 Enhanced Algorithm: Log Truncation, Checkpoints, Redo Optimization	473
13.3.4 The Complete Algorithm: Handling Transaction Aborts and Undo Completion	491
13.4 Redo-History Paradigm	501
13.4.1 Actions during Normal Operation	501
13.4.2 Simple Three-Pass and Two-Pass Algorithms	501
13.4.3 Enhanced Algorithms: Log Truncation, Checkpoints, and Redo Optimization	510
13.4.4 Complete Algorithms: Handling Transaction Rollbacks and Undo Completion	510
13.5 Lessons Learned	518
13.5.1 Putting Everything Together	519
Exercises	526
Bibliographic Notes	528
<b>Chapter 14</b> Object Model Crash Recovery	531
14.1 Goal and Overview	531
14.2 Conceptual Overview of Redo-History Algorithms	532
14.3 A Simple Redo-History Algorithm for Two-Layered Systems	536
14.3.1 Actions during Normal Operation	536
14.3.2 Steps during Restart	539
14.4 An Enhanced Redo-History Algorithm for Two-Layered Systems	545
14.5 A Complete Redo-History Algorithm for General Object Model Executions	552
14.6 Lessons Learned	556
Exercises	558
Bibliographic Notes	560

<b>Chapter 15</b>	<b>Special Issues of Recovery</b>	<b>561</b>
15.1	Goal and Overview	561
15.2	Logging and Recovery for Indexes and Large Objects	562
15.2.1	Logical Log Entries for the Redo of Index Page Splits	562
15.2.2	Logical Log Entries and Flush Ordering for Large-Object Operations	566
15.3	Intra-transaction Savepoints and Nested Transactions	571
15.4	Exploiting Parallelism during Restart	577
15.5	Special Considerations for Main-Memory Data Servers	580
15.6	Extensions for Data-Sharing Clusters	583
15.7	Lessons Learned	589
	Exercises	589
	Bibliographic Notes	591
<b>Chapter 16</b>	<b>Media Recovery</b>	<b>593</b>
16.1	Goal and Overview	593
16.2	Log-Based Method	596
16.2.1	Database Backup and Archive Logging during Normal Operation	597
16.2.2	Database Restore Algorithms	599
16.2.3	Analysis of the Mean Time to Data Loss	602
16.3	Storage Redundancy	606
16.3.1	Techniques Based on Mirroring	607
16.3.2	Techniques Based on Error-Correcting Codes	610
16.4	Disaster Recovery	618
16.5	Lessons Learned	620
	Exercises	621
	Bibliographic Notes	621
<b>Chapter 17</b>	<b>Application Recovery</b>	<b>623</b>
17.1	Goal and Overview	623
17.2	Stateless Applications Based on Queues	625
17.3	Stateful Applications Based on Queues	632
17.4	Workflows Based on Queues	637
17.4.1	Failure-Resilient Workflow State and Context	639
17.4.2	Decentralized Workflows Based on Queued Transactions	640
17.5	General Stateful Applications	642
17.5.1	Design Considerations	643
17.5.2	Overview of the Server Reply Logging Algorithm	646

17.5.3	Data Structures	648
17.5.4	Server Logging during Normal Operation	650
17.5.5	Client Logging during Normal Operation	653
17.5.6	Log Truncation	655
17.5.7	Server Restart	657
17.5.8	Client Restart	659
17.5.9	Correctness Reasoning	662
17.5.10	Applicability to Multi-tier Architectures	666
<b>17.6</b>	<b>Lessons Learned</b>	<b>667</b>
	Exercises	668
	Bibliographic Notes	669

## **PART FOUR**

### **COORDINATION OF DISTRIBUTED TRANSACTIONS**

<b>Chapter 18</b>	<b>Distributed Concurrency Control</b>	<b>673</b>
18.1	Goal and Overview	673
18.2	Concurrency Control in Homogeneous Federations	676
18.2.1	Preliminaries	676
18.2.2	Distributed 2PL	679
18.2.3	Distributed TO	680
18.2.4	Distributed SGT	683
18.2.5	Optimistic Protocols	685
18.3	Distributed Deadlock Detection	686
18.4	Serializability in Heterogeneous Federations	690
18.4.1	Global Histories	691
18.4.2	Global Serializability	694
18.4.3	Quasi Serializability	696
18.5	Achieving Global Serializability through Local Guarantees	698
18.5.1	Rigorousness	698
18.5.2	Commitment Ordering	700
18.6	Ticket-Based Concurrency Control	702
18.6.1	Explicit Tickets for Forcing Conflicts	702
18.6.2	Implicit Tickets	706
18.6.3	Mixing Explicit and Implicit Tickets	707
18.7	Object Model Concurrency Control in Heterogeneous Federations	708
18.8	Coherency and Concurrency Control for Data-Sharing Systems	710
18.9	Lessons Learned	716

Exercises	717
Bibliographic Notes	719
<b>Chapter 19 Distributed Transaction Recovery</b>	<b>723</b>
19.1 Goal and Overview	723
19.2 The Basic Two-Phase Commit Algorithm	725
19.2.1 2PC Protocol	725
19.2.2 Restart and Termination Protocol	733
19.2.3 Independent Recovery	741
19.3 The Transaction Tree Two-Phase Commit Algorithm	744
19.4 Optimized Algorithms for Distributed Commit	748
19.4.1 Presumed-Abort and Presumed-Commit Protocols	749
19.4.2 Read-Only Subtree Optimization	756
19.4.3 Coordinator Transfer	758
19.4.4 Reduced Blocking	761
19.5 Lessons Learned	763
Exercises	765
Bibliographic Notes	766
 <b>PART FIVE</b>	
<b>APPLICATIONS AND FUTURE PERSPECTIVES</b>	
<b>Chapter 20 What Is Next?</b>	<b>771</b>
20.1 Goal and Overview	771
20.2 What Has Been Achieved?	771
20.2.1 Ready-to-Use Solutions for Developers	772
20.2.2 State-of-the-Art Techniques for Advanced System Builders	773
20.2.3 Methodology and New Challenges for Researchers	775
20.3 Data Replication for Ubiquitous Access	776
20.4 E-Services and Workflows	779
20.5 Performance and Availability Guarantees	783
Bibliographic Notes	787
 <b>References</b>	<b>791</b>
<b>Index</b>	<b>829</b>
<b>About the Authors</b>	<b>853</b>

# Preface

Teamwork is essential. It allows you to blame someone else.

—Anonymous

## The Book's Mission

For three decades transaction processing has been a cornerstone of modern information technology: it is an indispensable asset in banking, stock trading, airlines, travel agencies, and so on. With the new millennium's proliferation of e-Commerce applications, business-to-business workflows, and broad forms of Web-based e-Services, transactional information systems are becoming even more important. Fortunately, the success of the transaction concept does not solely rely on clever system implementations, but builds on and leverages scientifically rigorous foundations that have been developed in the research community. This scientific achievement has most prominently been recognized by the 1998 Turing Award to Jim Gray for his outstanding, pioneering work on the transaction concept. It is exactly such a systematic and fundamental understanding that will allow us to generalize and extend transactional information systems toward the evolving new classes of network-centric, functionally rich applications.

For the above reason this book emphasizes scientific fundamentals of long-term validity and value, and does not cover specific system products, which tend to become quickly outdated. The book does, however, put the presented theory, algorithms, and implementation techniques into perspective with practical system architectures. In this sense, the book is complementary to the systems-oriented literature, most notably, the "TP bible" by Jim Gray and Andreas Reuter and the more recent textbook by Phil Bernstein and Eric Newcomer. Our role model instead is the classic book *Concurrency Control and Recovery in Database Systems*, by Phil Bernstein, Vassos Hadzilacos, and Nat Goodman, which is now out of print. However, the field has made much progress since the time that book was written, and the transaction concept has become of much broader relevance beyond the scope of database systems alone. Our book reflects the advances of the past decade and the trends in modern IT architectures.

## Organization of the Book

The two key components of a transactional information system are *concurrency control*, to ensure the correctness of data when many clients simultaneously access shared data, and *recovery*, to protect the data against system failures. The book devotes its two major parts, Part II and Part III, to these two components, organized into 15 chapters altogether. For distributed, multi-tier federations of transactional servers, we will show that the concurrency control and recovery components of each server are the major asset toward viable solutions, but in addition, the *coordination of distributed transactions* becomes a vital issue that will be covered in Part IV. These three technically “hard-core” parts are surrounded by Part I, which contains motivation and background material and outlines the “big picture” of transactional technology, and Part V, which gives an outlook on topics that could not be covered (for lack of space and time) and speculates on future trends. Throughout all five parts, each chapter begins with a brief section on the goal and overview of the chapter, and concludes with the sections Lessons Learned, Exercises, and Bibliographic Notes. (Note: Chapter 20 doesn’t include Lessons Learned or Exercises.)

## Guidelines for Teaching

This book covers advanced material, including intensive use of formal models. So a solid background in computer science in general is assumed, but not necessarily familiarity with database systems. Whatever knowledge from that area is needed will be provided within the book itself. It is, in fact, one of our major points that transactional technology is important for many other areas, such as operating systems, workflow management, electronic commerce, and distributed objects, and should therefore be taught independently of database classes.

The book is primarily intended as a text for advanced undergraduate courses or graduate courses, but we would also encourage industrial researchers as well as system architects and developers who need an in-depth understanding of transactional information systems to work with this book. After all, engineers should not be afraid of (a little bit of) mathematics.

The book has been class tested at both the University of the Saarland in Saarbrücken and the University of Münster, and partly also at the University of Constance, all in Germany, for advanced undergraduate courses. In Saarbrücken the course was organized in 15 teaching weeks, each with four hours lecturing and additional student assignments. A possible, approximate breakdown of the material for this teaching time frame is given below. Since many universities will allow only two hours of weekly lecturing for such an advanced course, the material can be divided into mandatory core subjects and optional “high-end” issues, as suggested in Table P.1 (with the first and last sections of each chapter always being mandatory and thus omitted in the table).

**Table P.1** Suggested teaching schedule for 15-week course.

<i>Week</i>	<i>Mandatory sections</i>	<i>Optional sections</i>
1	Chapter 1: <i>What Is It All About?</i> Chapter 2: <i>Computational Models</i> 1.2–1.5, 2.2–2.5	
2	Chapter 3: <i>Concurrency Control: Notions of Correctness for the Page Model</i> 3.2–3.5, 3.7–3.8	3.6, 3.9–3.10
3	Chapter 4: <i>Concurrency Control Algorithms</i> 4.2, 4.3.1–4.3.4, 4.4–4.5	4.3.5–4.3.8
4	Chapter 5: <i>Multiversion Concurrency Control</i> 5.2, 5.3.1–5.3.2, 5.5	5.3.3, 5.4
5	Chapter 6: <i>Concurrency Control on Objects: Notions of Correctness</i> 6.2–6.5	6.6
6	Chapter 7: <i>Concurrency Control Algorithms on Objects</i> Chapter 8: <i>Concurrency Control on Relational Databases</i> 7.2–7.5	7.6, 8.2–8.4
7	Chapter 9: <i>Concurrency Control on Search Structures</i> Chapter 10: <i>Implementation and Pragmatic Issues</i> 9.2–9.3, 9.4.1, 10.2–10.3	9.4.2, 9.4.3, 9.5, 10.4–10.7
8	Chapter 11: <i>Transaction Recovery</i> 11.2–11.4, 11.5.1, 11.6–11.7	11.5.2–11.5.3
9	Chapter 12: <i>Crash Recovery: Notion of Correctness</i> Chapter 13: <i>Page Model Crash Recovery Algorithms</i> 12.2–12.4, 13.2, 13.3.1–13.3.3, 13.4	12.5, 13.3.4
10	Chapter 14: <i>Object Model Crash Recovery</i> 14.2–14.4	14.5
11	Chapter 15: <i>Special Issues of Recovery</i> Chapter 16: <i>Media Recovery</i> 16.2.1–16.2.2, 16.3.1	15.2–15.6, 16.2.3, 16.3.2, 16.4
12	Chapter 17: <i>Application Recovery</i> 17.2–17.4	17.5
13	Chapter 18: <i>Distributed Concurrency Control</i> 18.2–18.3, 18.4.1–18.4.2, 18.5	18.4.3, 18.6–18.8
14	Chapter 19: <i>Distributed Transaction Recovery</i> 19.2–19.3	19.4
15	Chapter 20: <i>What Is Next?</i> 20.2–20.5	



**Table P.2** Suggested teaching schedule for 10-week course.

<i>Week</i>	<i>Mandatory sections</i>
1	Chapter 1: <i>What Is It All About?</i> Chapter 2: <i>Computational Models</i> 1.2–1.4, 2.2–2.5
2	Chapter 3: <i>Concurrency Control: Notions of Correctness for the Page Model</i> 3.2–3.5, 3.7–3.8
3	Chapter 4: <i>Concurrency Control Algorithms</i> 4.2, 4.3.1–4.3.4, 4.4–4.5
4	Chapter 5: <i>Multiversion Concurrency Control</i> 5.2, 5.3.1–5.3.2, 5.5
5	Chapter 6: <i>Concurrency Control on Objects: Notions of Correctness</i> Chapter 7: <i>Concurrency Control Algorithms on Objects</i> 6.2–6.5, 7.2–7.4
6	Chapter 10: <i>Implementation and Pragmatic Issues</i> Chapter 11: <i>Transaction Recovery</i> 10.2–10.3, 11.2–11.4, 11.5.1, 11.6–11.7
7	Chapter 12: <i>Crash Recovery: Notion of Correctness</i> Chapter 13: <i>Page Model Crash Recovery Algorithms</i> 12.2–12.4, 13.2, 13.3.1–13.3.3, 13.4
8	Chapter 14: <i>Object Model Crash Recovery</i> Chapter 16: <i>Media Recovery</i> 14.2–14.4, 16.2.1–16.2.2, 16.3.1
9	Chapter 17: <i>Application Recovery</i> Chapter 18: <i>Distributed Concurrency Control</i> 17.2–17.3, 18.2.1–18.2.2, 18.4.1–18.4.2, 18.5
10	Chapter 19: <i>Distributed Transaction Recovery</i> 19.2–19.3

It is also feasible to configure a 10-week course from the book's material. Under such time constraints it is obviously necessary to leave out some of the most advanced topics. Our subjective recommendations for a 10-week course, with either four or two hours lecturing per week, are shown in Table P.2.

Additional teaching materials, most notably, slides for lecturers and solutions to selected exercises are available at [www.mkp.com/tis/](http://www.mkp.com/tis/). We will also offer errata of the book as we discover our errors. And we'd appreciate comments, suggestions, and criticisms via email at [weikum@cs.uni-sb.de](mailto:weikum@cs.uni-sb.de) or [vossen@uni-muenster.de](mailto:vossen@uni-muenster.de).

## **Acknowledgments**

A number of colleagues provided us with very valuable input: encouragement, constructive criticism, proofreading, and class testing, and also simple “bug fixes.” We are most grateful to Jim Gray, Elliot Moss, Dennis Shasha, Betty and Pat O’Neil, K. Vidyasankar, Alan Fekete, Dave Lomet, and Marc Scholl. Significant work on the book’s exercises was contributed by Ralf Schenkel; Carolin Letz helped in editing and tracing the numerous bibliographic entries. The book has also implicitly benefited from many technical discussions and collaborations with Catriel Beeri, Yuri Breitbart, Theo Härder, Dave Lomet, and, most notably, Hans-Jörg Schek. Needless to say, all biases and possible errors in this book are our own.

Our editor Diane Cerra and her colleague Belinda Breyer were perfect in their balance between keeping us relaxed and creative and occasionally putting some healthy pressure on us, and they were always responsive to our needs. We wish everybody who is writing a book such a great editorial team. We were also lucky to cooperate with an excellent production team headed by Edward Wade, and including our copyeditor Judith Brown, proofreader Jennifer McClain, designers Rebecca Evans and Frances Baca, and indexer Steve Rath. Last but not least we would like to thank our families for being with us while we were mentally somewhere else.

This Page Intentionally Left Blank

## PART ONE

# **Background and Motivation**

This Page Intentionally Left Blank

# What Is It All About?

If I had had more time, I could have written you a shorter letter.

—Blaise Pascal

There are two mistakes one can make along the road  
to truth—not going all the way, and not starting.

—Buddha

### 1.1 Goal and Overview

Transaction processing is an important topic in database and information systems. Moreover, it is rapidly gaining importance outside the context in which it was originally developed. In this introductory chapter, we discuss why transactions are a good idea, why transactions form a reasonable abstraction concept for certain classes of real-life data management and related problems, as well as what can and what cannot be done with the transaction concept.

The transaction concept was originally developed in the context of database management systems as a paradigm for dealing with concurrent accesses to a shared database and for handling failures. Therefore, we start out (in Section 1.2) by describing typical application scenarios for database and other information systems in which transactions make sense. The original and most canonical application example is funds transfer in banking; very similar applications in terms of functionality and structure have arisen in a number of other service-providing industries, most notably in the travel industry with its flight, car, and hotel bookings. All these classical application examples are commonly referred to as *online transaction processing*, or *OLTP* for short. In addition, we will show that the application area of the transaction concept includes modern business sectors such as electronic commerce and the management of workflows (which are also known as business processes).

*Application  
areas*

In terms of the underlying computer and network infrastructure, we are typically dealing with distributed systems of potentially large scale and with possibly heterogeneous, interoperating components. Most often, one of these components is a database management system or, more specifically, a *database server* that processes requests issued by clients (workstations, personal

computers, portable notebooks, PDAs, electronic sensors, and other embedded systems). It turns out that in today's diverse information technology landscapes, mail servers, Web- or intranet-based document servers, and workflow management systems also play an increasingly important role and call for transactional support.

*Transaction  
concept*

The key problem that the transaction concept solves in a very elegant way is to cope with the subtle and often difficult issues of keeping data consistent even in the presence of highly concurrent data accesses and despite all sorts of failures. An additional key property of transactions is that this is achieved in a generic way that is essentially invisible to the application logic (and to application development), so that application developers are completely freed from the burden of dealing with such system issues. This is why transactions are an *abstraction concept*, and why this concept is a cornerstone of modern information technology. Section 1.3 will discuss the role of the transaction concept in state-of-the-art information systems from a strategic viewpoint. We will introduce a fairly general reference architecture as a bird's-eye view of the entire infrastructure that is necessary to implement and deploy an information system, and we will discuss several variations of this reference architecture that are commonly used in practice. In particular, we will identify components that are in charge of managing persistent data under a transaction-oriented access regime, and we will concentrate on these *transactional (data) servers*. We will then discuss, in Section 1.4, the abstract properties that constitute the transaction concept and the great benefit that these properties provide in the context of a transactional data server. We will also outline the requirements on the server's algorithms in terms of *correctness* and *performance*, as well as of *reliability* and *availability*.

*Computational  
models*

By far the most important concrete instantiation of a transactional data server is a database system. However, this is not a book about database systems. We limit our discussion to topics that are directly and closely related to transactions, and nothing else. We will briefly survey the kind of knowledge we expect our readers to have about database systems in Section 1.5. This will prepare the setting for the introduction of two computational models for transactional servers in the next chapter.

This chapter, like all subsequent chapters, is wrapped up by summarizing, in Section 1.6, the key insights that the reader should have obtained from reading it.

## 1.2 Application Examples

We begin our exposition with a few examples of applications in which transactional properties can be brought to bear; these scenarios are

funds transfer in a banking environment, a classical OLTP application, Web-based electronic commerce (e-Commerce), travel planning as a workflow example.

*From OLTP to  
e-Commerce  
and workflow*

### 1.2.1 Online Transaction Processing: Debit/Credit Example

Consider the simplified operation of a bank that uses a relational database for keeping track of its account business. The database contains, among others, a table named Account that describes bank accounts in terms of their account ID, associated customer name, identification of the respective bank branch, and balance. Transactions in the bank are either *withdrawals* or *deposits* (which is why the application is often characterized as consisting of debit/credit transactions), and these transactions are often combined in *funds transfers*. The typical structure of a debit/credit program is shown below, using commands of the standardized database query language SQL and embedding these commands in a C program. Note the distinction between local variables of the invoked program and the data in the underlying database that is shared by all programs. Also note that a realistic, full-fledged debit/credit program may include various sanity checks against the account data (e.g., for high amounts of withdrawals) between the SQL Select command and the subsequent Update step.

```
/* debit/credit program */
void main()
{
    EXEC SQL BEGIN DECLARE SECTION;
    int accountid, amount; /* input variables */
    int balance; /* intermediate variable */
    EXEC SQL END DECLARE SECTION;
    /* read user input */
    printf("Enter Account ID, Amount
    for deposit (positive) or withdrawal (negative):");
    scanf("%d%d", &accountid, &amount);
    /* determine current balance of the account,
    reading it into a local, intermediate, variable of the
    program */
    EXEC SQL Select Account.Balance Into :balance
    From Account
    Where Account.Id = :accountid;
    /* add amount (negative for withdrawal) */
    balance = balance + amount;
    /* update account balance in the database */
}
```



```

EXEC SQL Update Account
  Set Account_Balance = balance
  Where Account_Id = :accountid;
EXEC SQL Commit Work;
}

```

Assume the bank operates in a traditional way, with various tellers at its local branches executing transactions. However, the role of a teller's terminal could also be replaced by a customer's PC equipped with home banking software, the credit card reader of a merchant, or some other form of smart cybercash carrier. We will generally refer to these as "clients" and will disregard the details of the client software, as these are irrelevant to the server on which the database resides.

*Concurrency  
and parallelism  
→ concurrency  
control  
techniques for  
isolation*

With a huge number of clients potentially issuing simultaneous requests to the bank's database server, *concurrent* (i.e., overlapping in time) or even *parallel* (i.e., on multiple processors) execution of multiple debit/credit transactions is mandatory in order to exploit the server's hardware resources. For example, while the server is waiting for the completion of a disk I/O on behalf of one transaction, its CPU should be utilized to process another transaction; similarly, multiple transactions should be processed in parallel on a multiprocessor machine. Thus, the "outside world" of an individual transaction is continuously changing as concurrently executing transactions are modifying the underlying database of the banking application. In order to be able to ignore the potential fallacies of this concurrency, it is therefore desirable that each transaction be executed in an *isolated* manner, that is, as if there were no other transactions and hence no concurrency. We will show that this tension between concurrency for the sake of performance, on the one hand, and potential sequential execution for the sake of simplicity and correctness, on the other, is reconciled by the *concurrency control techniques* of a transactional server.

The following scenario illustrates that concurrency is indeed trickier than it may seem at first glance, and that it may have a disastrous impact on the consistency of the underlying data and thus the quality of the entire information system, even if each individual transaction is perfectly correct and preserves data consistency.

---

### EXAMPLE 1.1

Consider two debit/credit transactions (i.e., invocations of the debit/credit program outlined above) that are concurrently executed by processes  $P_1$  and  $P_2$ , respectively, both operating on the same account  $x$  (i.e., the value of the program's input variable "accountid"). To distinguish the two different instances of the local program variable "balance" that temporarily

holds the value of the account balance, we refer to them as *balance1* for process  $P_1$  and *balance2* for  $P_2$ . For simplicity, we ignore some syntactic details of the embedded SQL commands. The first transaction intends to withdraw \$30, and the second transaction intends to deposit \$20. We assume that the initial account balance is \$100. The table below shows those parts of the two transactions that read and modify the account record.

$P_1$	Time	$P_2$
/* balance1 = 0, x.Account_Balance = 100, balance2 = 0 */		
Select Account_Balance		
Into :balance1 From Account	1	
Where Account_Id = x		
/* balance1 = 100, x.Account_Balance = 100, balance2 = 0 */		
	2	Select Account_Balance
		Into :balance2 From Account
		Where Account_Id = x
/* balance1 = 100, x.Account_Balance = 100, balance2 = 100 */		
balance1 = balance1 - 30	3	
/* balance1 = 70, x.Account_Balance = 100, balance2 = 100 */		
	4	balance2 = balance2 + 20
/* balance1 = 70, x.Account_Balance = 100, balance2 = 120 */		
Update Account		
Set Account_Balance = :balance	5	
Where Account_Id = x		
/* balance1 = 70, x.Account_Balance = 70, balance2 = 120 */		
	6	Update Account
		Set Account_Balance = :balance2
		Where Account_Id = x
/* balance1 = 70, x.Account_Balance = 120, balance2 = 120 */		

Upon completion of the execution, the balance of account  $x$ , as recorded in the persistent database, will be \$120, although it should be \$90 after execution of the two transactions. Thus, the recorded data no longer reflects reality and should be considered incorrect. Obviously, for such an information system to be meaningful and practically viable, this kind of anomaly must be prevented by all means. Thus, concurrent executions must be treated with extreme care. Similar anomalies could arise from failures of processes or entire computers during the execution of a transaction, and need to be addressed as well.

A second fundamentally important point is that the various accesses that a transaction has to perform need to occur *in conjunction*. In other words,

*Failures*  
→ *recovery*  
*techniques for*  
*atomicity and*  
*durability*

once a transaction has started, its data accesses should look to the outside world as an atomic operation that is either executed completely or not at all. This property of *atomicity* will turn out to be a crucial requirement on database transactions. Moreover, this conceptual property should be guaranteed to hold even in a failure-prone environment where individual processes or the entire database server may fail at an arbitrarily inconvenient point in time. To this end, a transactional server provides *recovery techniques* to cope with failures. In addition to ensuring transaction atomicity, these techniques serve to ensure the *durability* of a transaction's effects once the transaction is completed.

The following scenario illustrates that atomicity is a crucial requirement for being able to cope with failures.

---

### EXAMPLE 1.2

Consider the following funds transfer program, which transfers a given amount of money between two accounts, by first withdrawing it from a source account and then depositing it in a target account. The program is described in terms of SQL statements embedded into a host program written in C.

```
/* funds transfer program */
void main()
{
    EXEC SQL BEGIN DECLARE SECTION;
        int sourceid, targetid, amount; /* input variables */
    EXEC SQL END DECLARE SECTION;
/* read user input */
    printf("Enter Source ID, Target ID, Amount to be
        transferred:");
    scanf("%d %d %d", &sourceid, &targetid, &amount);
/* subtract desired amount from source */
    EXEC SQL Update Account
        Set Account.Balance = Account.Balance - :amount
        Where Account.Id = :sourceid;
/* add desired amount to target */
    EXEC SQL Update Account
        Set Account.Balance = Account.Balance + :amount
        Where Account.Id = :targetid;
    EXEC SQL Commit Work;
}
```

Now assume that the above funds transfer program has started executing and has already performed its first update statement, withdrawing the

specified amount of money from the source. If there is a computer hardware or software failure that interrupts the program's execution at this critical point, the remaining second part will not be performed anymore. Thus, the target account will not receive the money, so that money is effectively lost in transit.

A recovery procedure, to be invoked after the system is restarted, could try to find out which updates were already made by ongoing transaction program executions and which ones were not yet done, and could try to fix the situation in some way. However, implementing such recovery procedures on a per-application-case basis is an extremely difficult task that is itself error prone by its mere complexity, especially because multiple transactions issued by different programs may have accessed the data at the time of the failure. So rather than programming recovery in an ad hoc manner for each application separately, a systematic approach is needed.

System-provided recovery that ensures the atomicity of transactions greatly simplifies the understanding of the postfailure state of the data and the overall failure handling on the application side. In the example scenario, rather than being left with the inconsistent state in the middle of the transaction, the system recovery should restore the state as of before the transaction began. On the other hand, if the transaction had already issued its "commit transaction" call and had received a positive returncode from the system, then the "all" case of the all-or-nothing paradigm would apply, and the system would henceforth guarantee the durability of the transaction's complete funds transfer.

---

The above conceptual properties of a transaction—namely, atomicity, durability, and isolation—together provide the key abstraction that allows application developers to disregard concurrency and failures, yet the transactional server guarantees the consistency of the underlying data and ultimately the correctness of the application. In the banking example, this means that no money is ever lost in the jungle of electronic funds transfers and customers can perfectly rely on electronic receipts, balance statements, and so on. As we will show in the next two application scenarios, these cornerstones for building highly dependable information systems can be successfully applied outside the scope of OLTP and classical database applications as well.

### 1.2.2 Electronic Commerce Example

In today's information landscape, client requests may span multiple databases and other information sources across enterprise boundaries, yet the mutual consistency of all this data is crucial and thus important to maintain. Then,

the resulting transactions operate in a distributed system that consists of multiple servers, often with heterogeneous software. As a concrete example of such a modern setting, consider what happens when a client intends to purchase something from an Internet-based bookstore; such applications are known as *electronic commerce* (e-Commerce).

The purchasing activity proceeds in the following steps:

1. The client connects to the bookstore's server through an appropriate Internet protocol, and starts browsing and querying the store's catalog.
2. The client gradually fills an electronic shopping cart with items that she intends to purchase.
3. When the client is about to check out, she reconsiders the items in her shopping cart and makes a final decision on which items she will purchase.
4. The client provides all the necessary information for placing a definitive (and legally binding) order. This includes her shipping address and information on her credit card or some other valid form of cybercash. The latter information may be encrypted such that the merchant can only verify its authenticity, but possibly without being able to actually decrypt the provided data.
5. The merchant's server forwards the payment information to the customer's bank, credit card company, or some other clearinghouse for cybercash. When the payment is accepted by the clearinghouse, the shipping of the ordered items is initiated by the merchant's server, and the client is notified on the successful completion of the e-Commerce activity.

So why are transactions and their properties relevant for this scenario? It is obviously important to keep certain data consistent, and this data is even distributed across different computers. The consistency requirement is already relevant during the catalog browsing phase when the user fills her shopping cart, as the client's view of the shopping cart should ideally be kept consistent with the shopping cart contents as maintained by the merchant's server. Note that this should be satisfied in the presence of temporary failures at the client or the server side (e.g., a software failure of a client or server process) and also network failures (e.g., undelivered messages due to network congestion). Further note that this seemingly simple requirement may transitively involve additional data, say, on the inventory for the selected items, which could reside on yet another computer.

While it could be argued that data consistency is merely an optional luxury feature for the shopping cart contents and does not necessarily justify the use of advanced technology like transactions in the technical sense of this book, a very similar situation arises in the last step of the entire activity. There, it is

absolutely crucial that three parties agree on the data that tracks the overall outcome:

The merchant's server must have records on both the order and the successfully certified payment.

At the same time, the clearinghouse must have a record on the payment, as its approval may be requested again later, or the clearinghouse may be responsible for the actual money transfer.

Finally, the client must have received the notification that the ordered items are being shipped.

When these three effects on three different computers are known to be atomic, confidence in the correct processing of such e-Commerce activities is greatly increased. Conversely, when atomicity is not guaranteed, all sorts of complicated cases arise, such as the merchant shipping the items but the clearinghouse losing all records of its cybercash approval and ultimately not being able to reclaim the money. Similarly, when the customer is never informed about the shipping and the resulting money transfer, she may order the items again from a different merchant, ending up with two copies of the same book. Even worse, the customer may receive the shipped items and keep them, but pretend that she never ordered and never received them. Then, it is obviously important for the entire e-Commerce industry to rely on atomicity guarantees in order to prove a customer's order when it comes to a lawsuit.

Similar, yet more involved arguments can be brought up about isolation properties, but the case for transactions should have been made sufficiently clear at this point. Of course, we could deal with inconsistent data among the three computers of our scenario in many other ways as well. But the decisive point is that by implementing the last step of the activity as a transaction, all the arguments about atomicity in the presence of failures can be factored out, and the entire application is greatly simplified.

So this example has indeed much more in common with the debit/credit scenario than it might have seemed at first glance. There are, however, a number of important differences as well, and these nicely highlight the potential generalization of the transaction concept beyond the classical setting of centralized database applications:

The entire application is *distributed* across multiple computers, and the software may be *heterogeneous* in that different database systems are used at the various servers. (Of course, the hardware is likely to be heterogeneous, too, but this is mostly masked by the software and thus less relevant.)

The servers are not necessarily based on database systems; they may as well be some other form of *information repository* or *document management servers* in general.

The effects of a transaction may even include *messages* between computers, for example, the notification of the customer. So transactions are not limited to what is usually perceived as “stored data.”

We will show in this book that transaction technology can cope well with all these additional, challenging aspects of modern applications. It will take us to some of the advanced material, however, to cover all issues.

### 1.2.3 Workflow Management: Travel Planning Example

A final and most challenging application class that we consider is so-called workflows, also known as (the computerized part of) *business processes*. A *workflow* is a set of activities (or steps) that belong together in order to achieve a certain business goal. Typical examples would be the processing of a credit request or insurance claim in a bank or insurance company, respectively; the work of a program committee for a scientific conference (submissions, reviews, notifications, etc.); the administrative procedures for real estate purchase; or the “routing” of a patient in a hospital. To orchestrate such processes, it is crucial to specify (at least a template for) the control flow and the data flow between activities, although it may still be necessary to improvise at run time (e.g., in medical applications). Making the “flow of work” between activities explicit in that it is factored out of the entire application is exactly the key leverage from workflow technology that allows a company or other institution to largely automate the repetitive, stereotypic parts of its processes while retaining flexibility, and to quickly adjust these processes to changing business needs.

Activities can be completely automated or based on interaction with a human user and intellectual decision making. This implies that workflows can be long lived, up to several days or weeks, or even months and years. A typical characteristic of workflows is that the activities are distributed across different responsible persons and different, independent information systems, possibly across different enterprises. In particular, an activity can spawn requests to an arbitrary “invoked application” that is provided by some server independently of the current workflow. Thus, workflow management is essentially an umbrella for the activities and invoked applications that constitute a particular workflow. To this end, a workflow management system provides a specification environment for registering activities and for specifying, in a high-level declarative way, not only the control and data flow within a process, but also a run-time environment that automatically triggers activities according to the specified flow. Workflow management systems with such capabilities are commercially available and are gaining significant industrial relevance.

As a concrete example of a workflow, consider the activities that are necessary in the planning of a business trip, say, a trip to a conference. Suppose your manager (or professor) allows you to choose one scientific or developer’s

conference that you can attend as part of a continuing-education program. This involves the following activities:

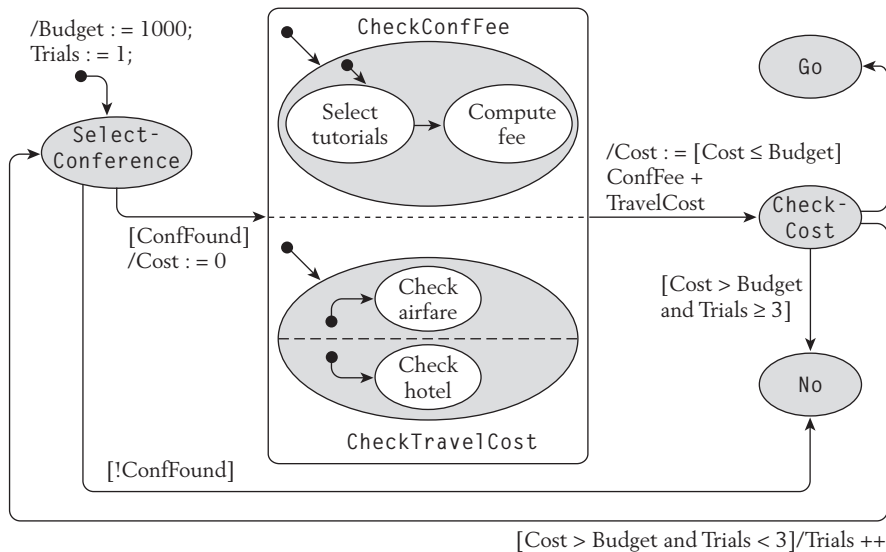
Select a conference, based on its subject, technical program, time, and place. If no suitable conference is found, the process is terminated.

Check out the cost of the trip to this conference, typically by delegation to a travel agency.

Check out the registration fee for the conference, which often depends on your memberships, tutorials that you may wish to attend, and so on.

Compare the total cost of attending the selected conference to the allowed budget, and decide to attend the conference only if the cost is within the budget.

With the increasing costs of conferences and ever tighter travel budgets (at the time this book was written), it is desirable to allow several trials with different conferences, but the number of trials should be limited, in order to guarantee termination of the entire process. The activities and the control flow between them are graphically depicted in Figure 1.1. This illustration is based on a specification formalism known as *statecharts*, which is one particular kind of formal specification method that might be used by a workflow management system. Each oval denotes a state in which the workflow can exist during its execution. Each state in turn corresponds to one activity; so



**Figure 1.1** Specification of the travel planning workflow.



the activity is spawned when the state is entered. The activity may then invoke further application programs. When the workflow is started, a specified initial state (a state without predecessors) is entered, and the workflow terminates when a final state (a state without successors) is reached. In the example, the initial state is the `SelectConference` state, and the final states are `Go` and `No`.

The transitions between states are governed by event-condition-action rules that are attached to the transition arcs as labels. The meaning of a rule of the form  $E[C]/A$  is that the transition fires if event  $E$  has occurred and condition  $C$  is true in the current state. Then the current state is left and the state where the transition arc points to is entered; during this transition the specified action  $A$  is performed. In the example, we only make use of conditions and actions. Both refer to a small set of variables (instantiated for each workflow instance) that are relevant for the control flow. This kind of control flow specification allows conditional execution as well as loops based on high-level predicates. The entire specification can be hierarchical, thus supporting both top-down refinement and bottom-up composition of existing building blocks, by allowing states to be nested. So a state can in turn contain another statechart that is executed when the state is entered. In addition, the specification formalism allows parallel execution, which is graphically indicated by breaking a state down into two or more *orthogonal* statecharts, separated by a dashed line, that are executed in parallel. In the example, the activities that correspond to the two states `CheckConfFee` and `CheckTravelCost` are executed in parallel. These two states are further refined into several steps, where `CheckTravelCost` again leads to two parallel substates.

Although the example scenario is still largely oversimplified, the above discussion already indicates some of the semantically rich process design issues that accompany workflow management. Here we are interested in the connection between workflows and transactions, and how a workflow application could possibly benefit from transaction-supporting services. The answer is threefold and involves different stages of transactional scope:

The activities themselves can, of course, spawn requests to information systems that lead to transactional executions in these systems. This is almost surely the case with the `CheckTravelCost` activity. The travel agency's invoked application would typically issue transactions against the reservation systems of airlines and hotels. In fact, it seems to make sense that this activity not only figures out the prices, but also makes reservations in the underlying information systems. Obviously, booking a flight to a certain city and a hotel room in that city makes sense only if both reservations are successful. If either of the two is unavailable, the whole trip no longer makes sense. So these two steps need to be tied together in a single transaction. Note that this transaction is a distributed one that involves two autonomous information systems.

The outcome of the above reservations affects the further processing of the workflow. The requests against the various information systems would return status codes that should be stored in variables of the workflow and would be relevant for the future control flow. For example, not being able to make one of the two necessary reservations in the selected city should trigger going back to the initial `SelectConference` state for another trial. (To keep the example specification simple, this is not shown in Figure 1.1.) Thus, it is desirable (if not mandatory) that the modification of the workflow's variables be embedded in the same transaction that accesses the airline and hotel databases. In other words, the *state of the workflow application* should be under transactional control as well. This is an entirely new aspect that did not arise in the banking and e-Commerce examples. It is questionable whether today's commercial workflow management systems can cope with this issue in the outlined, transactional way. But as we will show, transactional technology does provide solutions for incorporating application state into atomic processing units as well.

We could discuss whether the entire travel planning workflow should be a single transaction that incorporates all effects on the underlying information systems as well as the state of the workflow application itself. After all, the entire workflow should have an all-or-nothing, atomic effect. Ideas along these lines have indeed been discussed in the research community for quite a few years; however, no breakthrough is in sight. The difficulty lies in the long-lived nature of workflows and the fact that workflows, like simple transactions, run concurrently. Atomicity is therefore coupled with isolation properties: the atomicity of a workflow would imply that no concurrent workflow could ever "see" any partial effects of it. Regardless of the technical details of how isolation can be implemented at all (to be covered in great depth in this book), maintaining such isolation over a period of hours, days, or weeks raises questions about performance problems with regard to the progress of concurrent workflows. For this reason, the straightforward approach of turning an entire workflow into a single transaction is absolutely infeasible.

The discussion of the third item above does not imply, however, that the one-transaction-per-activity approach is the only kind of transactional support for workflows. Consider the situation when all necessary reservations have been made successfully, but later it is found that the total cost including the conference fees is unacceptable and it is decided not to attend any conference at all. Now you hold reservations that may later result in charges to your credit card, unless you intervene. So you must make sure that these reservations are canceled. One approach could be to extend the workflow specification by additional cancellation activities and the necessary control flow. However, it turns out that cancellation-like activities are fairly common in many business

processes, and a lot of time would be spent in specifying these kinds of things over and over again. So a better solution would be to generalize the particular case at hand into a more abstract notion of *compensation activities*. Each compensation activity would be tied to one of the regular activities of the workflow in the sense that their combined effect is a neutral one from the workflow application's viewpoint. Of course, it would still be necessary to provide code (or a high-level declarative specification) for each compensation activity itself, but modifications to the workflow's control flow specification are no longer required. Instead, the appropriate triggering of compensation activities could be delegated to the workflow management system. The transactional technology that we develop in this book does provide the principal means for coping with compensation issues in the outlined, generic way (as opposed to developing specific solutions on a per-application basis over and over again).

At this point in the book, the major insight from this discussion is to realize that the scope of transactions is not a priori fixed and limited to stored data, but can be (carefully) extended to incorporate various aspects of information system applications as well.

## 1.3 System Paradigms

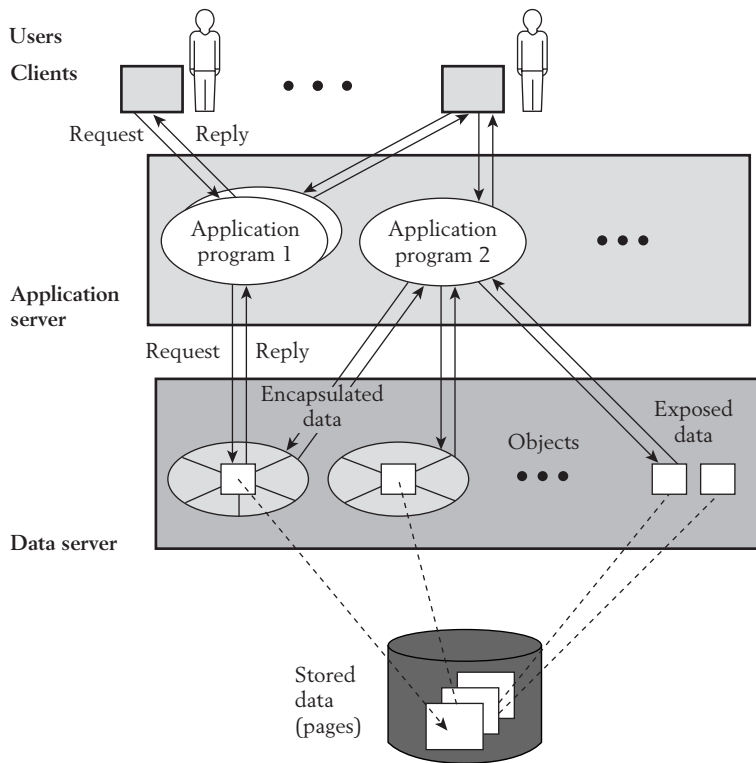
The three application scenarios of Section 1.2 already provide hints on the various system components that it takes to implement and deploy such applications in their entirety and on putting them together into a complete system architecture. In particular, we have seen that we need to separate *clients*—that is, the computers or terminals from which a human user generates computer work—from the *servers* where data and possibly executable programs reside in various forms. However, this distinction alone is insufficient for characterizing full-fledged modern information systems. We now introduce a more systematic view of these architectural issues in that we set up a *reference architecture* (or framework) to capture the most typical cases that are used in practice.

### 1.3.1 Three-Tier and Two-Tier Architectures

*Reference  
architecture:  
three-tier  
system*

Our reference architecture is illustrated in Figure 1.2. It captures what is most frequently referred to by practitioners as a *three-tier architecture*. It consists of a set of clients (PCs, workstations, notebooks, terminals, digital TV set top boxes, “intelligent” sensors, etc.) that interact with an *application server*, which in turn interacts with a *data server*.

Clients send business-oriented (or goal-oriented) requests to the application server. For example, invoking a debit/credit transaction, starting an



**Figure 1.2** Reference architecture.

e-Commerce shopping session, or initiating a travel planning workflow would be concrete examples for such requests. In modern applications, requests are typically issued from a GUI (Graphical User Interface); likewise, the reply that the application server will send back is often presented in a graphical way, using forms, buttons, charts, or even virtual reality-style animations. All this presentation processing, for both input and output, is done by the client. Therefore, HTML (Hypertext Markup Language), one of the original cornerstones of the World Wide Web, is a particularly attractive basis for presentation, because it merely requires that a Web browser is installed on the client side and thus applies to a large set of clients.

The application server has a repository of *application programs* in executable form, and invokes the proper program that is capable of handling a client request. Both the application programs and the entire application server can be organized in a variety of ways. The programs themselves may be anything from an old-fashioned terminal-oriented COBOL program to a Java applet or some other program that generates, for example, a dynamic page.

*Clients:*  
*presentation*  
*processing*

*Application*  
*server*