

Physical Database Design

The Morgan Kaufmann Series in Data Management Systems

Series Editor: Jim Gray, Microsoft Research

Data Preparation for Data Mining Using SAS
Mamdouh Refaat

Querying XML: XQuery, XPath, and SQL/XML in Context

Jim Melton and Stephen Buxton

Data Mining: Concepts and Techniques,
Second Edition

Jiawei Han and Micheline Kamber

Database Modeling and Design: Logical Design,
Fourth Edition

Toby J. Teorey, Sam S. Lightstone and Thomas P. Nadeau

Foundations of Multidimensional and Metric Data Structures

Hanan Samet

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

Moving Objects Databases

Ralf Hartmut Güting and Markus Schneider

Joe Celko's SQL Programming Style
Joe Celko

Data Mining, Second Edition: Concepts and Techniques

Ian Witten and Eibe Frank

Fuzzy Modeling and Genetic Algorithms for Data Mining and Exploration
Earl Cox

Data Modeling Essentials, Third Edition
Graeme C. Simsion and Graham C. Witt

Location-Based Services

Jochen Schiller and Agnès Voisard

Database Modeling with Microsoft® Visio for Enterprise Architects

Terry Halpin, Ken Evans, Patrick Hallock, Bill Maclean

Designing Data-Intensive Web Applications
Stephano Ceri, Piero Fraternali, Aldo Bongio, Marco Brambilla, Sara Comai, and Maristella Matera

Mining the Web: Discovering Knowledge from Hypertext Data

Soumen Chakrabarti

Advanced SQL: 1999—Understanding Object-Relational and Other Advanced Features
Jim Melton

Database Tuning: Principles, Experiments, and Troubleshooting Techniques

Dennis Shasha and Philippe Bonnet

SQL:1999—Understanding Relational Language Components

Jim Melton and Alan R. Simon

Information Visualization in Data Mining and Knowledge Discovery

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

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

Gerhard Weikum and Gottfried Vossen

Spatial Databases: With Application to GIS

Philippe Rigaux, Michel Scholl, and Agnes Voisard

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

Terry Halpin

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

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 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

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

Philip A. Bernstein and Eric Newcomer

Using the New DB2: IBMs 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 O₂

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

Physical Database Design

The Database Professional's Guide
to Exploiting Indexes, Views,
Storage, and More

Sam Lightstone
Toby Teorey
Tom Nadeau



AMSTERDAM • BOSTON • HEIDELBERG • LONDON
NEW YORK • OXFORD • PARIS • SAN DIEGO
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO
Morgan Kaufmann Publishers is an imprint of Elsevier



MORGAN KAUFMANN PUBLISHERS

| | |
|------------------------------------|--------------------------|
| <i>Publisher</i> | Diane D. Cerra |
| <i>Publishing Services Manager</i> | George Morrison |
| <i>Project Manager</i> | Marilyn E. Rash |
| <i>Assistant Editor</i> | Asma Palmeiro |
| <i>Cover Image</i> | Nordic Photos |
| <i>Composition:</i> | Multiscience Press, Inc. |
| <i>Interior Printer</i> | Sheridan Books |
| <i>Cover Printer</i> | Phoenix Color Corp. |

Morgan Kaufmann Publishers is an imprint of Elsevier.
500 Sansome Street, Suite 400, San Francisco, CA 94111

 This book is printed on acid-free paper.

Copyright © 2007 by Elsevier Inc. All rights reserved.

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.

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, scanning, or otherwise—without prior written permission of the publisher.

Permissions may be sought directly from Elsevier's Science & Technology Rights Department in Oxford, UK: phone: (+44) 1865 843830, fax: (+44) 1865 853333, E-mail: permissions@elsevier.com. You may also complete your request on-line via the Elsevier homepage (<http://elsevier.com>), by selecting "Support & Contact" then "Copyright and Permission" and then "Obtaining Permissions."

Page 197: "Make a new plan Stan . . . and get yourself free." - Paul Simon, copyright (c) Sony BMG/Columbia Records. All rights reserved. Used with permission.

Library of Congress Cataloging-in-Publication Data

Lightstone, Sam.

Physical database design : the database professional's guide to exploiting indexes, views, storage, and more / Sam Lightstone, Toby Teorey, and Tom Nadeau.

p. cm. -- (The Morgan Kaufmann series in database management systems)

Includes bibliographical references and index.

ISBN-13: 978-0-12-369389-1 (alk. paper)

ISBN-10: 0-12-369389-6 (alk. paper)

1. Database design. I. Teorey, Toby J. II. Nadeau, Tom, 1958– III. Title.

QA76.9.D26L54 2007

005.74--dc22

2006102899

For information on all Morgan Kaufmann publications, visit our Web site at www.mkp.com or www.books.elsevier.com

Printed in the United States of America

07 08 09 10 11 10 9 8 7 6 5 4 3 2 1

Working together to grow
libraries in developing countries

www.elsevier.com | www.bookaid.org | www.sabre.org

ELSEVIER

BOOK AID
International

Sabre Foundation

Contents

| | |
|---|--------------|
| Preface | xv |
| Organization | xvi |
| Usage Examples | xvii |
| Literature Summaries and Bibliography | xviii |
| Feedback and Errata | xviii |
| Acknowledgments | xix |
| I Introduction to Physical Database Design | I |
| I.1 Motivation—The Growth of Data and Increasing Relevance of Physical Database Design | 2 |
| I.2 Database Life Cycle | 5 |
| I.3 Elements of Physical Design: Indexing, Partitioning, and Clustering | 7 |
| I.3.1 Indexes | 8 |
| I.3.2 Materialized Views | 9 |
| I.3.3 Partitioning and Multidimensional Clustering | 10 |
| I.3.4 Other Methods for Physical Database Design | 10 |
| I.4 Why Physical Design Is Hard | 11 |
| I.5 Literature Summary | 12 |

| | | |
|----------|---|-----------|
| 2 | Basic Indexing Methods | 15 |
| 2.1 | B+tree Index | 16 |
| 2.2 | Composite Index Search | 20 |
| 2.2.1 | Composite Index Approach | 24 |
| 2.2.2 | Table Scan | 24 |
| 2.3 | Bitmap Indexing | 25 |
| 2.4 | Record Identifiers | 27 |
| 2.5 | Summary | 28 |
| 2.6 | Literature Summary | 28 |
| 3 | Query Optimization and Plan Selection | 31 |
| 3.1 | Query Processing and Optimization | 32 |
| 3.2 | Useful Optimization Features in Database Systems | 32 |
| 3.2.1 | Query Transformation or Rewrite | 32 |
| 3.2.2 | Query Execution Plan Viewing | 33 |
| 3.2.3 | Histograms | 33 |
| 3.2.4 | Query Execution Plan Hints | 33 |
| 3.2.5 | Optimization Depth | 34 |
| 3.3 | Query Cost Evaluation—An Example | 34 |
| 3.3.1 | Example Query 3.1 | 34 |
| 3.4 | Query Execution Plan Development | 41 |
| 3.4.1 | Transformation Rules for Query Execution Plans | 42 |
| 3.4.2 | Query Execution Plan Restructuring Algorithm | 42 |
| 3.5 | Selectivity Factors, Table Size, and Query Cost Estimation | 43 |
| 3.5.1 | Estimating Selectivity Factor for a Selection Operation or Predicate | 43 |
| 3.5.2 | Histograms | 45 |
| 3.5.3 | Estimating the Selectivity Factor for a Join | 46 |
| 3.5.4 | Example Query 3.2 | 46 |
| 3.5.5 | Example Estimations of Query Execution Plan Table Sizes | 49 |
| 3.6 | Summary | 50 |
| 3.7 | Literature Summary | 51 |
| 4 | Selecting Indexes | 53 |
| 4.1 | Indexing Concepts and Terminology | 53 |
| 4.1.1 | Basic Types of Indexes | 54 |
| 4.1.2 | Access Methods for Indexes | 55 |
| 4.2 | Indexing Rules of Thumb | 55 |

| | | |
|----------|--|-----------|
| 4.3 | Index Selection Decisions | 58 |
| 4.4 | Join Index Selection | 62 |
| 4.4.1 | Nested-loop Join | 62 |
| 4.4.2 | Block Nested-loop Join | 65 |
| 4.4.3 | Indexed Nested-loop Join | 65 |
| 4.4.4 | Sort-merge Join | 66 |
| 4.4.5 | Hash Join | 67 |
| 4.5 | Summary | 69 |
| 4.6 | Literature Summary | 70 |
| 5 | Selecting Materialized Views | 71 |
| 5.1 | Simple View Materialization | 72 |
| 5.2 | Exploiting Commonality | 77 |
| 5.3 | Exploiting Grouping and Generalization | 84 |
| 5.4 | Resource Considerations | 86 |
| 5.5 | Examples: The Good, the Bad, and the Ugly | 89 |
| 5.6 | Usage Syntax and Examples | 92 |
| 5.7 | Summary | 95 |
| 5.8 | Literature Review | 96 |
| 6 | Shared-nothing Partitioning | 97 |
| 6.1 | Understanding Shared-nothing Partitioning | 98 |
| 6.1.1 | Shared-nothing Architecture | 98 |
| 6.1.2 | Why Shared Nothing Scales So Well | 100 |
| 6.2 | More Key Concepts and Terms | 101 |
| 6.3 | Hash Partitioning | 101 |
| 6.4 | Pros and Cons of Shared Nothing | 103 |
| 6.5 | Use in OLTP Systems | 106 |
| 6.6 | Design Challenges: Skew and Join Collocation | 108 |
| 6.6.1 | Data Skew | 108 |
| 6.6.2 | Collocation | 109 |
| 6.7 | Database Design Tips for Reducing Cross-node Data Shipping | 110 |
| 6.7.1 | Careful Partitioning | 110 |
| 6.7.2 | Materialized View Replication and Other Duplication Techniques | 111 |
| 6.7.3 | The Internode Interconnect | 115 |
| 6.8 | Topology Design | 117 |
| 6.8.1 | Using Subsets of Nodes | 117 |
| 6.8.2 | Logical Nodes versus Physical Nodes | 119 |

| | | |
|------|----------------------|-----|
| 6.9 | Where the Money Goes | 120 |
| 6.10 | Grid Computing | 120 |
| 6.11 | Summary | 121 |
| 6.12 | Literature Summary | 122 |

7 Range Partitioning 125

| | | |
|-------|--|-----|
| 7.1 | Range Partitioning Basics | 126 |
| 7.2 | List Partitioning | 128 |
| 7.2.1 | Essentials of List Partitioning | 128 |
| 7.2.2 | Composite Range and List Partitioning | 128 |
| 7.3 | Syntax Examples | 129 |
| 7.4 | Administration and Fast Roll-in and Roll-out | 131 |
| 7.4.1 | Utility Isolation | 131 |
| 7.4.2 | Roll-in and Roll-out | 133 |
| 7.5 | Increased Addressability | 134 |
| 7.6 | Partition Elimination | 135 |
| 7.7 | Indexing Range Partitioned Data | 138 |
| 7.8 | Range Partitioning and Clustering Indexes | 139 |
| 7.9 | The Full Gestalt: Composite Range and Hash Partitioning with Multidimensional Clustering | 139 |
| 7.10 | Summary | 142 |
| 7.11 | Literature Summary | 142 |

8 Multidimensional Clustering 143

| | | |
|-------|--|-----|
| 8.1 | Understanding MDC | 144 |
| 8.1.1 | Why Clustering Helps So Much | 144 |
| 8.1.2 | MDC | 145 |
| 8.1.3 | Syntax for Creating MDC Tables | 151 |
| 8.2 | Performance Benefits of MDC | 151 |
| 8.3 | Not Just Query Performance: Designing for Roll-in and Roll-out | 152 |
| 8.4 | Examples of Queries Benefiting from MDC | 153 |
| 8.5 | Storage Considerations | 157 |
| 8.6 | Designing MDC Tables | 159 |
| 8.6.1 | Constraining the Storage Expansion Using Coarsification | 159 |
| 8.6.2 | Monotonicity for MDC Exploitation | 162 |
| 8.6.3 | Picking the Right Dimensions | 163 |
| 8.7 | Summary | 165 |
| 8.8 | Literature Summary | 166 |

| | | |
|-----------|--|------------|
| 9 | The Interdependence Problem | 167 |
| 9.1 | Strong and Weak Dependency Analysis | 168 |
| 9.2 | Pain-first Waterfall Strategy | 170 |
| 9.3 | Impact-first Waterfall Strategy | 171 |
| 9.4 | Greedy Algorithm for Change Management | 172 |
| 9.5 | The Popular Strategy (the Chicken Soup Algorithm) | 173 |
| 9.6 | Summary | 175 |
| 9.7 | Literature Summary | 175 |
| 10 | Counting and Data Sampling in Physical Design Exploration | 177 |
| 10.1 | Application to Physical Database Design | 178 |
| 10.1.1 | Counting for Index Design | 180 |
| 10.1.2 | Counting for Materialized View Design | 180 |
| 10.1.3 | Counting for Multidimensional Clustering Design | 182 |
| 10.1.4 | Counting for Shared-nothing Partitioning Design | 183 |
| 10.2 | The Power of Sampling | 184 |
| 10.2.1 | The Benefits of Sampling with SQL | 184 |
| 10.2.2 | Sampling for Database Design | 185 |
| 10.2.3 | Types of Sampling | 189 |
| 10.2.4 | Repeatability with Sampling | 192 |
| 10.3 | An Obvious Limitation | 192 |
| 10.4 | Summary | 194 |
| 10.5 | Literature Summary | 195 |
| 11 | Query Execution Plans and Physical Design | 197 |
| 11.1 | Getting from Query Text to Result Set | 198 |
| 11.2 | What Do Query Execution Plans Look Like? | 201 |
| 11.3 | Nongraphical Explain | 201 |
| 11.4 | Exploring Query Execution Plans to Improve Database Design | 205 |
| 11.5 | Query Execution Plan Indicators for Improved Physical Database Designs | 211 |
| 11.6 | Exploring without Changing the Database | 214 |
| 11.7 | Forcing the Issue When the Query Optimizer Chooses Wrong | 215 |
| 11.7.1 | Three Essential Strategies | 215 |

| | | |
|--------|--|-----|
| 11.7.2 | Introduction to Query Hints | 216 |
| 11.7.3 | Query Hints When the SQL Is Not Available to Modify | 219 |
| 11.8 | Summary | 220 |
| 11.9 | Literature Summary | 220 |

12 Automated Physical Database Design 223

| | | |
|--------|---|-----|
| 12.1 | What-if Analysis, Indexes, and Beyond | 225 |
| 12.2 | Automated Design Features from Oracle, DB2, and SQL Server | 229 |
| 12.2.1 | IBM DB2 Design Advisor | 231 |
| 12.2.2 | Microsoft SQL Server Database Tuning Advisor | 234 |
| 12.2.3 | Oracle SQL Access Advisor | 238 |
| 12.3 | Data Sampling for Improved Statistics during Analysis | 240 |
| 12.4 | Scalability and Workload Compression | 242 |
| 12.5 | Design Exploration between Test and Production Systems | 247 |
| 12.6 | Experimental Results from Published Literature | 248 |
| 12.7 | Index Selection | 254 |
| 12.8 | Materialized View Selection | 254 |
| 12.9 | Multidimensional Clustering Selection | 256 |
| 12.10 | Shared-nothing Partitioning | 258 |
| 12.11 | Range Partitioning Design | 260 |
| 12.12 | Summary | 262 |
| 12.13 | Literature Summary | 262 |

13 Down to the Metal: Server Resources and Topology 265

| | | |
|--------|---|-----|
| 13.1 | What You Need to Know about CPU Architecture and Trends | 266 |
| 13.1.1 | CPU Performance | 266 |
| 13.1.2 | Amdahl's Law for System Speedup with Parallel Processing | 269 |
| 13.1.3 | Multicore CPUs | 271 |
| 13.2 | Client Server Architectures | 271 |
| 13.3 | Symmetric Multiprocessors and NUMA | 273 |
| 13.3.1 | Symmetric Multiprocessors and NUMA | 273 |
| 13.3.2 | Cache Coherence and False Sharing | 274 |
| 13.4 | Server Clusters | 275 |
| 13.5 | A Little about Operating Systems | 275 |

| | | |
|-----------|--|------------|
| 13.6 | Storage Systems | 276 |
| 13.6.1 | Disks, Spindles, and Striping | 277 |
| 13.6.2 | Storage Area Networks and Network Attached Storage | 278 |
| 13.7 | Making Storage Both Reliable and Fast Using RAID | 279 |
| 13.7.1 | History of RAID | 279 |
| 13.7.2 | RAID 0 | 281 |
| 13.7.3 | RAID 1 | 281 |
| 13.7.4 | RAID 2 and RAID 3 | 282 |
| 13.7.5 | RAID 4 | 284 |
| 13.7.6 | RAID 5 and RAID 6 | 284 |
| 13.7.7 | RAID 1+0 | 285 |
| 13.7.8 | RAID 0+1 | 285 |
| 13.7.9 | RAID 10+0 and RAID 5+0 | 286 |
| 13.7.10 | Which RAID Is Right for Your Database Requirements? | 288 |
| 13.8 | Balancing Resources in a Database Server | 288 |
| 13.9 | Strategies for Availability and Recovery | 290 |
| 13.10 | Main Memory and Database Tuning | 295 |
| 13.10.1 | Memory Tuning by Mere Mortals | 295 |
| 13.10.2 | Automated Memory Tuning | 298 |
| 13.10.3 | Cutting Edge: The Latest Strategy in Self-tuning Memory Management | 301 |
| 13.11 | Summary | 314 |
| 13.12 | Literature Summary | 314 |
| 14 | Physical Design for Decision Support, Warehousing, and OLAP | 317 |
| 14.1 | What Is OLAP? | 318 |
| 14.2 | Dimension Hierarchies | 320 |
| 14.3 | Star and Snowflake Schemas | 321 |
| 14.4 | Warehouses and Marts | 323 |
| 14.5 | Scaling Up the System | 327 |
| 14.6 | DSS, Warehousing, and OLAP Design Considerations | 328 |
| 14.7 | Usage Syntax and Examples for Major Database Servers | 329 |
| 14.7.1 | Oracle | 330 |
| 14.7.2 | Microsoft's Analysis Services | 331 |
| 14.8 | Summary | 333 |
| 14.9 | Literature Summary | 334 |

| | | |
|-------------------|---|------------|
| 15 | Denormalization | 337 |
| 15.1 | Basics of Normalization | 338 |
| 15.2 | Common Types of Denormalization | 342 |
| 15.2.1 | Two Entities in a One-to-One Relationship | 342 |
| 15.2.2 | Two Entities in a One-to-many Relationship | 343 |
| 15.3 | Table Denormalization Strategy | 346 |
| 15.4 | Example of Denormalization | 347 |
| 15.4.1 | Requirements Specification | 347 |
| 15.4.2 | Logical Design | 349 |
| 15.4.3 | Schema Refinement Using Denormalization | 350 |
| 15.5 | Summary | 354 |
| 15.6 | Literature Summary | 354 |
| 16 | Distributed Data Allocation | 357 |
| 16.1 | Introduction | 358 |
| 16.2 | Distributed Database Allocation | 360 |
| 16.3 | Replicated Data Allocation—"All-beneficial Sites" | |
| | Method | 362 |
| 16.3.1 | Example | 362 |
| 16.4 | Progressive Table Allocation Method | 367 |
| 16.5 | Summary | 368 |
| 16.6 | Literature Summary | 369 |
| Appendix A | A Simple Performance Model for Databases | 371 |
| A.1 | I/O Time Cost—Individual Block Access | 371 |
| A.2 | I/O Time Cost—Table Scans and Sorts | 372 |
| A.3 | Network Time Delays | 372 |
| A.4 | CPU Time Delays | 374 |
| Appendix B | Technical Comparison of DB2 HADR with Oracle Data Guard for Database Disaster Recovery | 375 |
| B.1 | Standby Remains "Hot" during Failover | 376 |
| B.2 | Subminute Failover | 377 |
| B.3 | Geographically Separated | 377 |
| B.4 | Support for Multiple Standby Servers | 377 |
| B.5 | Support for Read on the Standby Server | 377 |
| B.6 | Primary Can Be Easily Reintegrated after Failover | 378 |

| | |
|--------------------------|------------|
| Glossary | 379 |
| Bibliography | 391 |
| Index | 411 |
| About the Authors | 427 |

Preface

Since the development of the relational model by E. F. Codd at IBM in 1970, relational databases have become the de facto standard for managing and querying structured data. The rise of the Internet, online transaction processing, online banking, and the ability to connect heterogeneous systems have all contributed to the massive growth in data volumes over the past 15 years. Terabyte-sized databases have become commonplace. Concurrent with this data growth have come dramatic increases in CPU performance spurred by Moore's Law, and improvements in disk technology that have brought about a dramatic increase in data density for disk storage. Modern databases frequently need to support thousands if not tens of thousands of concurrent users. The performance and maintainability of database systems depends dramatically on their physical design.

A wealth of technologies has been developed by leading database vendors allowing for a fabulous range of physical design features and capabilities. Modern databases can now be sliced, diced, shuffled, and spun in a magnificent set of ways, both in memory and on disk. Until now, however, not much has been written on the topic of physical database design. While it is true that there have been white papers and articles about individual features and even individual products, relatively little has been written on the subject as a whole. Even less has been written to commiserate with database designers over the practical difficulties that the complexity of "creeping featurism" has imposed on the industry. This is all the more reason why a text on physical database design is urgently needed.

We've designed this new book with a broad audience in mind, with both students of database systems and industrial database professionals clearly within its scope. In it

we introduce the major concepts in physical database design, including indexes (B+, hash, bitmap), materialized views (deferred and immediate), range partitioning, hash partitioning, shared-nothing design, multidimensional clustering, server topologies, data distribution, underlying physical subsystems (NUMA, SMP, MPP, SAN, NAS, RAID devices), and much more. In keeping with our goal of writing a book that had appeal to students and database professionals alike, we have tried to concentrate the focus on practical issues and real-world solutions.

In every market segment and in every usage of relational database systems there seems to be nowhere that the problems of physical database design are not a critical concern: from online transaction processing (OLTP), to data mining (DM), to multidimensional online analytical processing (MOLAP), to enterprise resource planning (ERP), to management resource planning (MRP), and in both in-house enterprise systems designed and managed by teams of database administrators (DBAs) and in deployed independent software vendor applications (ISVAs). We hope that the focus on physical database design, usage examples, product-specific syntax, and best practice, will make this book a very useful addition to the database literature.

Organization

An overview of physical database design and where it fits into the database life cycle appears in Chapter 1. Chapter 2 presents the fundamentals of B+tree indexing, the most popular indexing method used in the database industry today. Both simple indexing and composite indexing variations are described, and simple performance measures are used to help compare the different approaches. Chapter 3 is devoted to the basics of query optimization and query execution plan selection from the viewpoint of what a database professional needs to know as background for database design.

Chapters 4 through 8 discuss the individual important design decisions needed for physical database design. Chapter 4 goes into the details about how index selection is done, and what alternative indexing strategies one has to choose from for both selection and join operations. Chapter 5 describes how one goes about choosing materialized views for individual relational databases as well as setting up star schemas for collections of databases in data warehouses. The tradeoffs involved in materialized view selection are illustrated with numerical examples. Chapter 6 explains how to do shared-nothing partitioning to divide and conquer large and computationally complex database problems. The relationship between shared-nothing partitioning, materialized view replication, and indexing is presented.

Chapter 7 is devoted to range partitioning, dividing a large table into multiple smaller tables that hold a specific range of data, and the special indexing problems that need to be addressed. Chapter 8 discusses the benefits of clustering data in general, and how powerful this technique can be when extended to multidimensional data. This

allows a system to cluster along multiple dimensions at the same time without duplicating data.

Chapter 9 discusses the problem of integrating the many physical design decisions by exploring how each decision affects the others, and leads the designer into ways to optimize the design over these many components. Chapter 10 looks carefully at methods of counting and sampling data that help improve the individual techniques of index design, materialized view selection, clustering, and partitioning. Chapter 11 goes more thoroughly into query execution plan selection by discussing tools that allow users to look at the query execution plans and observe whether database decisions on design choices, such as index selection and materialized views, are likely to be useful.

Chapter 12 contains a detailed description of how many of the important physical design decisions are automated by the major relational databases—DB2, SQL Server, and Oracle. It discusses how to use these tools to design efficient databases more quickly. Chapter 13 brings the database designer in touch with the many system issues they need to understand: multiprocessor servers, disk systems, network topologies, disaster recovery techniques, and memory management.

Chapter 14 discusses how physical design is needed to support data warehouses and the OLAP techniques for efficient retrieval of information from them. Chapter 15 defines what is meant by denormalization and illustrates the tradeoffs between degree of normalization and database performance. Finally, Chapter 16 looks at the basics of distributed data allocation strategies including the tradeoffs between the fast query response times due to data replication and the time cost of updates of multiple copies of data.

Appendix A briefly describes a simple computational performance model used to evaluate and compare different physical design strategies on individual databases. The model is used to clarify the tradeoff analysis and design decisions used in physical design methods in several chapters. Appendix B includes a comparison of two commercially available disaster-recovery technologies—IBM's High Availability Disaster Recovery and Oracle's Data Guard.

Each chapter has a tips and insights section for the database professional that gives the reader a useful summary of the design highlights of each chapter. This is followed by a literature summary for further investigation of selected topics on physical design by the reader.

Usage Examples

One of the major differences between logical and physical design is that with physical design the underlying features and physical attributes of the database server (its software and its hardware) begin to matter much more. While logical design can be performed in the abstract, somewhat independent of the products and components that will be used to materialize the design, the same cannot be said for physical design. For this reason we

have made a deliberate effort to include examples in this second book of the major database server products in database server products about physical database design. In this set we include DB2 for zOS v8.1, DB2 9 (Linux, Unix, and Windows), Oracle 10g, SQL Server 2005, Informix Dataserver, and NCR Teradata. We believe that this covers the vast majority of industrial databases in use today. Some popular databases are conspicuously absent, such as MySQL and Sybase, which were excluded simply to constrain the authoring effort.

Literature Summaries and Bibliography

Following the style of the our earlier text on logical database design, *Database Modeling and Design: Logical Design, Fourth Edition*, each chapter concludes with a literature summary. These summaries include the major papers and references for the material covered in the chapter, specifically in two forms:

- Seminal papers that represent the original breakthrough thinking for the physical database design concepts discussed in the chapter.
- Major papers on the latest research and breakthrough thinking.

In addition to the chapter-centric literature summaries, a larger more comprehensive bibliography is included at the back of this book.

Feedback and Errata

If you have comments, we would like to hear from you. In particular, it's very valuable for us to get feedback on both changes that would improve the book as well as errors in the current content. To make this possible we've created an e-mail address to dialogue with our readers: please write to us at db-design@rogers.com.

Has everyone noticed that all the letters of the word *database* are typed with the left hand? Now the layout of the QWERTY typewriter keyboard was designed among other things to facilitate the even use of both hands. It follows, therefore, that among other things, writing about databases is not only unnatural, but a lot harder than it appears.

—Anonymous

While this quip may appeal to the authors who had to personally suffer through left-hand-only typing of the word *database* several hundred times in the authoring of this book,¹ if you substitute the words “writing about databases” with “designing data-

¹ Confession of a bad typist: I use my right hand for the t and b. This is an unorthodox but necessary variation for people who need to type the word “database” dozens of times per day.

bases,” the statement rings even more powerfully true for the worldwide community of talented database designers.

Acknowledgments

As with any text of this breadth, there are many people aside from the authors who contribute to the reviewing, editing, and publishing that make the final text what it is. We’d like to pay special thanks to the following people from a range of companies and consulting firms who contributed to the book: Sanjay Agarwal, Eric Alton, Hermann Baer, Kevin Beck, Surajit Chaudhuri, Kitman Cheung, Leslie Cranston, Yuri Deigin, Chris Eaton, Scott Fadden, Lee Goddard, Peter Haas, Scott Hayes, Lilian Hobbs, John Hornibrook, Martin Hubel, John Kennedy, Eileen Lin, Guy Lohman, Wenbin Ma, Roman Melnyk, Mughees Minhas, Vivek Narasayya, Jack Raitto, Haider Rizvi, Peter Shum, Danny Zilio and Calisto Zuzarte. Thank you to Linda Peterson and Rebekah Smith for their help with manuscript preparation.

We also would like to thank the reviewers of this book who provided a number of extremely valuable insights. Their in-depth reviews and new directions helped us produce a much better text. Thank you to Mike Blaha, Philippe Bonnet, Philippe Carino, and Patrick O’Neil. Thank you as well to the concept reviewers Bob Muller, Dorian Pyle, James Bean, Jim Gray, and Michael Blaha.

We would like to thank our wives and children for their support and for allowing us the time to work on this project, often into the wee hours of the morning.

To the community of students and database designers worldwide, we salute you. Your job is far more challenging and complex than most people realize. Each of the possible design attributes in a modern relational database system is very complex in its own right. Tackling all of them, as real database designers must, is a remarkable challenge that by all accounts ought to be impossible for mortal human beings. In fact, optimal database design can be shown mathematically to truly be impossible for any moderately involved system. In one analysis we found that the possible design choices for an average database far exceeded the current estimates of the number of atoms in the universe (10^{81}) by several orders of magnitude! And yet, despite the massive complexity and sophistication of modern database systems, you have managed to study them, master them, and continue to design them. The world’s data is literally in your hands. We hope this book will be a valuable tool for you. By helping you, the students and designers of database systems, we hope this book will also lead in a small incremental but important way to improvements in the world’s data management infrastructure.

Engineering is a great profession. There is the satisfaction of watching a figment of the imagination emerge through the aid of science to a plan on paper. Then it moves to realization in stone or metal or energy. Then it brings homes to men or women. Then

it elevates the standard of living and adds to the comforts of life. This is the engineer's high privilege.

—*Herbert Hoover (1874–1964)*

The most likely way for the world to be destroyed, most experts agree, is by accident. That's where we come in; we're computer professionals. We cause accidents.

—*Nathaniel Borenstein (1957–)*