



**Mastering
Data Warehouse
Aggregates
Solutions for Star
Schema Performance**

Christopher Adamson



WILEY

Wiley Publishing, Inc.

Mastering Data Warehouse Aggregates



**Mastering
Data Warehouse
Aggregates
Solutions for Star
Schema Performance**

Christopher Adamson



WILEY

Wiley Publishing, Inc.

Mastering Data Warehouse Aggregates: Solutions for Star Schema Performance

Published by

Wiley Publishing, Inc.

10475 Crosspoint Boulevard

Indianapolis, IN 46256

www.wiley.com

Copyright © 2006 by Wiley Publishing, Inc., Indianapolis, Indiana

Published simultaneously in Canada

ISBN-13: 978-0-471-77709-0

ISBN-10: 0-471-77709-9

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1

1MA/SQ/QW/QW/IN

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, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. Requests to the Publisher for permission should be addressed to the Legal Department, Wiley Publishing, Inc., 10475 Crosspoint Blvd., Indianapolis, IN 46256, (317) 572-3447, fax (317) 572-4355, or online at <http://www.wiley.com/go/permissions>.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Website is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Website may provide or recommendations it may make. Further, readers should be aware that Internet Websites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services or to obtain technical support, please contact our Customer Care Department within the U.S. at (800) 762-2974, outside the U.S. at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Library of Congress Cataloging-in-Publication Data

Adamson, Christopher, 1967–

Mastering data warehouse aggregates: solutions for star schema performance / Christopher Adamson.

p. cm.

Includes index.

ISBN-13: 978-0-471-77709-0 (pbk.)

ISBN-10: 0-471-77709-9 (pbk.)

1. Data warehousing. I. Title.

QA76.9.D37A333 2006

005.74—dc22

2006011219

Trademarks: Wiley, the Wiley logo, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. All other trademarks are the property of their respective owners. Wiley Publishing, Inc., is not associated with any product or vendor mentioned in this book.

For Wayne H. Adamson

1929–2003

*Through those whose lives you touched,
your spirit of love endures.*



About the Author

Christopher Adamson is a data warehousing consultant and founder of Oakton Software LLC. An expert in star schema design, he has managed and executed data warehouse implementations in a variety of industries. His customers have included Fortune 500 companies, large and small businesses, government agencies, and data warehousing tool vendors. Mr. Adamson also teaches dimensional modeling and is a co-author of *Data Warehouse Design Solutions* (also from Wiley). He can be contacted through his website, www.ChrisAdamson.net.



Executive Editor

Robert Elliott

Development Editor

Brian Herrmann

Technical Editor

Jim Hadley

Copy Editor

Nancy Rapoport

Editorial Manager

Mary Beth Wakefield

Production Manager

Tim Tate

**Vice President and Executive
Group Publisher**

Richard Swadley

**Vice President and Executive
Publisher**

Joseph B. Wikert

Project Coordinator

Michael Kruzil

Graphics and Production

Specialists

Jennifer Click

Denny Hager

Stephanie D. Jumper

Heather Ryan

Quality Control Technicians

John Greenough

Brian H. Walls

Proofreading and Indexing

Techbooks



Contents

Foreword	xix
Acknowledgments	xxi
Introduction	xxiii
Chapter 1 Fundamentals of Aggregates	1
Star Schema Basics	2
Operational Systems and the Data Warehouse	3
Operational Systems	3
Data Warehouse Systems	4
Facts and Dimensions	5
The Star Schema	7
Dimension Tables and Surrogate Keys	7
Fact Tables and Grain	10
Using the Star Schema	13
Multiple Stars and Conformance	15
Data Warehouse Architecture	20
Invisible Aggregates	22
Improving Performance	23
The Base Schema and the Aggregate Schema	25
The Aggregate Navigator	26
Principles of Aggregation	27
Providing the Same Results	27
The Same Facts and Dimension Attributes as the Base Schema	28

Other Types of Summarization	29
Pre-Joined Aggregates	29
Derived Tables	30
Tables with New Facts	31
Summary	32
Chapter 2 Choosing Aggregates	35
What Is a Potential Aggregate?	36
Aggregate Fact Tables: A Question of Grain	36
Aggregate Dimensions Must Conform	37
Pre-Joined Aggregates Have Grain Too	39
Enumerating Potential Aggregates	39
Identifying Potentially Useful Aggregates	40
Drawing on Initial Design	41
Design Decisions	41
Listening to Users	44
Where Subject Areas Meet	45
The Conformance Bus	45
Aggregates for Drilling Across	46
Query Patterns of an Existing System	49
Analyzing Reports for Potential Aggregates	49
Choosing Which Reports to Analyze	54
Assessing the Value of Potential Aggregates	55
Number of Aggregates	55
Presence of an Aggregate Navigator	55
Space Consumed by Aggregate Tables	56
How Many Rows Are Summarized	57
Examining the Number of Rows Summarized	59
The Cardinality Trap and Sparsity	62
Who Will Benefit from the Aggregate	64
Summary	65
Chapter 3 Designing Aggregates	67
The Base Schema	68
Identification of Grain	68
When Grain Is Forgotten	68
Grain and Aggregates	69
Conformance Bus	70
Rollup Dimensions	72
Aggregation Points	74
Natural Keys	74
Source Mapping	75
Slow Change Processing	76
Hierarchies	76
Housekeeping Columns	78

Design Principles for the Aggregate Schema	81
A Separate Star for Each Aggregation	81
Single Schema and the Level Field	81
Drawbacks to the Single Schema Approach	84
Advantages of Separate Tables	85
Pre-Joined Aggregates	86
Naming Conventions	87
Naming the Attributes	87
Naming Aggregate Tables	88
Aggregate Dimension Design	90
Attributes of Aggregate Dimensions	90
Sourcing Aggregate Dimensions	91
Shared Dimensions	92
Aggregate Fact Table Design	93
Aggregate Facts: Names and Data Types	94
No New Facts, Including Counts	94
Degenerate Dimensions	96
Audit Dimension	96
Sourcing Aggregate Fact Tables	97
Pre-Joined Aggregate Design	98
Documenting the Aggregate Schema	98
Identify Schema Families	99
Identify Dimensional Conformance	99
Documenting Aggregate Dimension Tables	101
Documenting Aggregate Fact Tables	103
Pre-Joined Aggregates	106
Materialized Views and Materialized Query Tables	108
Summary	108
Chapter 4 Using Aggregates	109
Which Tables to Use?	110
The Schema Design	110
Relative Size	113
Aggregate Portfolio and Availability	114
Requirements for the Aggregate Navigator	116
Why an Aggregate Navigator?	116
Two Views and Query Rewrite	117
Dynamic Availability	120
Multiple Front Ends	121
Multiple Back Ends	123
Evaluating Aggregate Navigators	126
Front-End Aggregate Navigators	127
Approach	127
Pros and Cons	128

Back-End Aggregate Navigation	129
Approach	129
Pros and Cons	130
Performance Add-On Technologies and OLAP	134
Approach	134
Pros and Cons	135
Specific Solutions	136
Living with Materialized Views	137
Using Materialized Views	137
Materialized Views as Pre-Joined Aggregates	137
Materialized Views as Aggregate Fact Tables (Without Aggregate Dimensions)	140
Materialized Views and Aggregate Dimension Tables Additional Considerations	141
Living with Materialized Query Tables	144
Using Materialized Query Tables	144
Materialized Query Tables as Pre-Joined Aggregates	145
Materialized Query Tables as Aggregate Fact Tables (Without Aggregate Dimensions)	146
Materialized Query Tables and Aggregate Dimension Tables Additional Considerations	147
Working Without an Aggregate Navigator	148
Human Decisions	149
Maintaining the Aggregate Portfolio	150
Impact on the ETL Process	151
Summary	151
Chapter 5 ETL Part 1: Incorporating Aggregates	153
The Load Process	154
The Importance of the Load	154
Tools of the Load	155
Incremental Loads and Changed Data Identification	156
The Top-Level Process	157
Loading the Base Star Schema	158
Loading Dimension Tables	159
Attributes of the Dimension Table	159
Requirements for the Dimension Load Process	161
Extracting and Preparing the Record	161
Process New Records	163
Process Type 1 Changes	164
Process Type 2 Changes	165
Loading Fact Tables	167
Requirements for the Fact Table Load Process	167
Acquire Data and Assemble Facts	168
Identification of Surrogate Keys	170
Putting It All Together	172

Loading the Aggregate Schema	174
Loading Aggregates Separately from Base Schema Tables	174
Invalid Aggregates	175
Load Frequency	176
Taking Aggregates Off-Line	176
Off-Line Load Processes	177
Materialized Views and Materialized Query Tables	178
Drop and Rebuild Versus Incremental Load	180
Drop and Rebuild	180
Incremental Loading of Aggregates	181
Real-Time Loads	182
Real-Time Load of the Base Schema	182
Real-Time Load and Aggregate Tables	183
Partitioning the Schema	183
Summary	185
Chapter 6 ETL Part 2: Loading Aggregates	187
The Source Data for Aggregate Tables	188
Changed Data Identification	188
Elimination of Redundant Processing	189
Ensuring Conformance	190
Loading the Base Schema and Aggregates Simultaneously	192
Loading Aggregate Dimensions	193
Requirements for the Aggregate Dimension Load Process	194
Extracting and Preparing the Records	195
Identifying and Processing New Records	197
Identifying and Processing Type 1 Changes	198
Processing Type 2 Changes	203
Key Mapping	204
Loading Aggregate Fact Tables	205
Requirements for Loading Aggregate Fact Tables	205
Acquire Data and Assemble Facts	205
Selecting Source Columns	206
Processing New Facts Only	208
Calculating and Aggregating Facts	208
One Query Does It All	209
Identification of Surrogate Keys	210
Aggregating Over Time	212
Dropping and Rebuilding Aggregates	214
Dropping and Rebuilding Aggregate Dimension Tables	214
Dropping and Rebuilding Aggregate Fact Tables	216
Pre-Joined Aggregates	217
Dropping and Rebuilding a Pre-Joined Aggregate	217
Incrementally Loading a Pre-Joined Aggregate	219
Materialized Views and Materialized Query Tables	221
Defining Attributes for Aggregate Dimensions	221
Optimizing the Hierarchy	222
Summary	223

Chapter 7	Aggregates and Your Project	225
	Data Warehouse Implementation	226
	Incremental Implementation of the Data Warehouse	226
	Planning Data Marts Around Conformed Dimensions	226
	Other Approaches	230
	Incorporating Aggregates into the Project	230
	Aggregates and the First Data Mart	231
	Subsequent Subject Areas	232
	The Aggregate Project	233
	Strategy Stage	234
	Technology Selection: Choosing an Aggregate Navigator	234
	Additional Strategic Tasks and Deliverables	240
	Design Stage	241
	Design of the Aggregate Schema and Load Specification	243
	Design Documentation	243
	Developing Test Plans for Aggregates	244
	Build Stage	245
	Iterative Build and Aggregates	245
	Build Tasks and Aggregates	246
	Deployment	250
	Transitioning to Production, Final Testing, and Documentation	250
	End User Education	252
	Management of Aggregates	252
	Maintenance Responsibilities	252
	Ad Hoc Changes to Aggregate Portfolio	254
	An Ongoing Process	254
	Summary	255
Chapter 8	Advanced Aggregate Design	257
	Aggregating Facts	258
	Periodic Snapshots Designs	258
	Transactions	258
	Snapshots	259
	Semi-Additivity	260
	Invisible Aggregates for Periodic Snapshots	261
	Averaging Semi-Additive Facts Produces a Derived Schema	263
	Taking Less Frequent Snapshots Does Not Produce an Invisible Aggregate	264
	Accumulating Snapshots	265
	The Accumulating Snapshot	265
	Aggregating the Accumulating Snapshot	267
	Factless Fact Tables	269
	Factless Events and Aggregates	269
	Coverage Tables and Aggregates	271

Aggregating Dimensions	272
Transaction Dimensions	273
Timestamping a Dimension	273
Aggregating a Timestamped Dimension	274
Bridge Tables	275
Dealing with Multi-Valued Attributes	276
Aggregates and Bridge Tables	278
Core and Custom Stars	282
Other Schema Types	283
Snowflakes and Aggregates	283
The Snowflake Schema	283
Aggregating Snowflakes	284
Third Normal Form Schemas and Aggregates	287
Summary	290
Chapter 9 Related Topics	291
Aggregates and the Archive Strategy	292
The Data Warehouse Archive Strategy	292
Aggregates and Archives	295
Maintaining Aggregates	295
Archive Versus Purge	297
Summarizing Off-Line Data	297
Aggregates and Security	299
Dimensionally Driven Security and Aggregates	299
Unrestricted Access to Summary Data	301
Derived Tables	302
The Merged Fact Table	303
The Pivoted Fact Table	307
The Sliced Fact Table	309
When Rollups Are Deployed Before Detail	310
Building the Base Table First	311
Building the Rollup First	312
Parallel Load Processes	312
Redeveloping the Load	314
Historic Detail	316
Summary	317
Glossary	319
Index	329



Foreword

In 1998 I wrote the foreword for Chris Adamson and Mike Venerable's book *Data Warehouse Design Solutions* (Wiley, 1998). Over the intervening eight years I have been delighted to track that book, as it has stayed high in the list of data warehouse best sellers, even through today. Chris and Mike had identified a set of data warehouse design challenges and were able to speak very effectively in that book to the community of data warehouse designers.

Viewed in the right perspective, the mission of data warehousing has not changed at all since 1998! In that foreword, I wrote that the data warehouse must be driven from business analysis needs, must be a mirror of management's urgent priorities, and must be a presentation facility that is understandable and fast. All of these perspectives have held true through today. While our databases have exploded in size, and the database content has become much more operational, the original description of the data warehouse rings true. If anything, the data warehouse, in its role as the platform for all forms of business intelligence, has become much more important than it was in 1998.

At the same time that the reach of the data warehouse has penetrated to every worker's desktop, we have all been swept along by the development of the Internet, and particularly search engines like Google. This parallel revolution, surprisingly, has sent data warehousing and business intelligence a powerful and simple message. As the saying goes, "The medium is the message." In this case, Google's message is:

You can search the entire contents of the Internet in less than a second.

The message to data warehousing is:

You should expect instantaneous results from your data warehouse queries.

To be perfectly frank, data warehousing and business intelligence have so far made only partial progress toward instantaneous performance. Our databases are more complicated than Google's documents, and our queries are more complex. *But*, we have some powerful tools that can be used to get us much closer to the goal of instantaneous performance.

Those of us who, like Chris and the Kimball Group, have long recognized that the class of data warehouse designs known as dimensional models offers a systematic opportunity for a huge performance boost, above and beyond database indexes, hardware RAM, faster processors, or parallelism. In fact, this additional performance opportunity, known as *aggregates*, when used correctly, can trump all the other performance techniques!

The idea behind aggregates is very simple. Always start with the most atomic, transaction-grained data available from the original source systems. Place that atomic data in full view of the end users in a dimensional format. Of course, if you stop there, you will have performance problems because many queries will do a huge amount of I/O no matter how much hardware you throw at the problem. Now aggregates come to the rescue. You systematically create a set of physically stored, pre-calculated summary records that are predictable common queries, or parts of queries posed by the end users. These summary records are the aggregates.

Aggregates, when used correctly, can provide performance improvements of a hundred or even a thousand times. No other technology is capable of such gains.

This book is all about aggregates. Chris explains how they rely on the dimensional approach, which aggregates to build, how to build them, and how to maintain them. He also shows in detail how Oracle's materialized views and IBM's materialized query tables are perfect examples of aggregates used effectively.

I was delighted to see Chris return to being an author after his wonderful first book. His only excuse for waiting eight years was that he was "busy building data warehouses." I'll accept that excuse! Now we can apply Chris's insights into making our data warehouse and business intelligence systems a big step closer to being instantaneous.

Ralph Kimball
Founder, Kimball Group
Boulder Creek, California



Acknowledgments

Thank you to everyone who read my first book, *Data Warehouse Design Solutions*, which I wrote with Mike Venerable. The positive feedback we received from around the world was unexpected, and most appreciated. Without your warm reception, I doubt that the current volume would have come to pass.

This book would not have been possible without Ralph Kimball. The value of his contribution to the data warehousing world cannot be understated. He has established a practical and powerful approach to data warehousing and provided terminology and principles for dimensional modeling that are used throughout the industry. I am deeply grateful for Ralph's continued support and encouragement, without which neither this nor my previous book would have been written.

I thank everyone at Wiley who contributed to this effort. Bob Elliott was a pleasure to work with and provided constructive criticism that was instrumental in shaping this book. Brian Herrmann made the writing process as painless as possible. I also thank the anonymous reviewers of my original outline, whose comments made this a better book.

Thanks also to Jim Hadley, who put in long hours reviewing drafts of this book. Through his detailed comments and advice, he made a substantial contribution to this effort. His continuing encouragement got me through several rough spots.

I am grateful to the customers and colleagues with whom I have worked over the years. The opportunity to learn from one another enriches us all. In particular, I thank three people as yet unmentioned. Mike Venerable has offered me opportunities that have shaped my career, along with guidance and advice that have helped me grow in numerous dimensions. Greg Jones's

xxii Acknowledgments

work managing data warehouse projects has profoundly influenced my own perspective, as is evident in Chapter 7. And Randall Porter has always been a welcome source of professional guidance, which was offered over many breakfasts during the writing of this book.

A very special thank you to my wife, Gladys, and sons, Justin and Carter, whose support and encouragement gave me the resolve I needed to complete this project. I also received support from my mother, sister, in-laws, and sisters-in-law. I could not have done this without all of you.



Introduction

In the battle to improve data warehouse performance, no weapon is more powerful and efficient than the aggregate table. A well-planned set of aggregates can have an extraordinary impact on the overall throughput of the data warehouse. After you ensure that the database is properly designed, configured, and tuned, any measures taken to address data warehouse performance should begin with aggregates.

Yet many businesses continue to ignore aggregates, turning instead to proprietary hardware products, converting to specialized databases, or implementing complex caching architectures. These solutions carry high price tags for acquisition and implementation and often require specialized skills to maintain. This book aims to fill the knowledge gap that has led businesses down this expensive and risky path.

In these pages, you will find tools and techniques you can use to bring stunning performance gains to *your* data warehouse. This book develops a set of best practices for the selection, design, construction, and use of aggregate tables. It explores how these techniques can be incorporated into projects, studies advanced design considerations, and covers how aggregates affect other aspects of the data warehouse lifecycle.

Intended Audience

This book is intended for *you*, the data warehouse practitioner with an interest in improving query performance. You may serve any one of a number of roles in the data warehouse environment, including:

- Business analyst
- Star schema designer
- Report developer
- ETL developer
- Project manager
- Database administrator
- Data administrator
- I.T. director
- Chief information officer
- Power user

Regardless of your role or current level of knowledge, the best practices in this book will help you and your team achieve astounding increases in data warehouse performance, without requiring an investment in exotic technologies.

It will be assumed that you have a very basic familiarity with relational database technology, understanding the concepts of tables, columns, and joins. Occasional examples of SQL code will be provided, and they will be fully explained in the accompanying text.

For those new to data warehousing, the background necessary to understand the examples will be provided along the way. For example, an overview of the star schema is presented in Chapter 1. The Extract Transform Load (ETL) process for the data warehouse is described in Chapter 5. The high-level data mart implementation process is described in Chapter 7.

About This Book

This book assumes a star schema approach to data warehousing. The necessary background is provided for readers new to this approach. It also considers implications of snowflake designs and, to a lesser extent, schemas in third normal form (3NF).

The design principles and best practices developed in each chapter make no assumptions about specific software products in the data warehouse. This tool-agnostic perspective is periodically supplemented with specific advice for users of Oracle's materialized views and IBM DB/2's materialized query tables.

Star Schema Approach

The techniques presented in this book are intended for data warehouses that are designed according to the principles of *dimensional modeling*, more popularly known as the *star schema approach*. Popularized by Ralph Kimball in the 1990s, the dimensional model is now widely accepted as the optimal method to organize information for analytic consumption.

Ralph Kimball and Margy Ross provide a comprehensive treatment of dimensional modeling in *The Data Warehouse Toolkit, Second Edition* (Wiley, 2002). The seminal work on the subject, their book is required reading for any student of data warehousing. The best practices in this book build on the foundation provided by Kimball and Ross and are described using terminology established by *The Toolkit*.

If you are not familiar with the star schema approach to data warehouse design, Chapter 1 provides an overview of the basic principles necessary to understand the examples in this book.

Snowflakes and 3NF Designs

Although this book focuses on the star schema, it does not ignore other approaches to schema design. From time to time, this book will examine the impact of a *snowflake design* on principles established throughout the book. For example, implications of a snowflake schema for aggregate design are explored in Chapters 2 and 3, and discussed more fully in Chapter 8.

In addition, Chapter 8 will look at how dimensional aggregates can service a *third normal form* schema design. Because of the complex relationships between the tables of a normalized schema, dimensional aggregates can have a tremendous impact. Of course, this is really the impact of the dimensional model itself. Best practices would suggest beginning with the most granular design possible, which is not really an aggregate at all. Still, a dimensional perspective can be used to augment query performance in such an environment.

Tool Independence

This book makes no assumptions regarding the presence of specific software products in your data warehouse architecture. Many commercial products offer features to assist in the implementation of aggregate tables. Each implementation is different; each has its own benefits and drawbacks; all are constantly changing.

Regardless of the tools used to build and navigate aggregates, you will need to address the same major tasks. You must choose which aggregates to implement; the aggregates must be designed; the aggregates must be built; a process must be established to ensure they are refreshed, or loaded, on a regular basis; the warehouse must be configured so that application queries are redirected to the aggregates.

This book provides a set of principles and best practices to guide you through these common tasks.

You can also use the principles in this book to guide the selection of specific technologies. For example, one component that you may need to add to your data warehouse architecture is the *aggregate navigator*. Chapter 4 develops a set of requirements for the aggregate navigator function. Three styles of commercial implementations are identified and evaluated against these requirements. You can use these requirements to evaluate your current technology options, as described in Chapter 7. They will remain valid even as specific products change and evolve.

Materialized Views and Materialized Query Tables

Specific database features from Oracle (materialized views) and IBM's DB/2 (materialized query tables) can be used to load and maintain aggregate tables as well as provide aggregate navigation services.

Throughout this book, the impact of using these technologies to build and navigate dimensional aggregates is explored. After establishing principles and best practices, we consider the implications of using these products. What is potentially gained or lost? How can you modify your process to accommodate the products' strengths and weaknesses? This is information that cannot be gleaned from a syntax reference manual.

Keep in mind that these products continue to evolve. Over time, their capabilities can be expected to expand and change. If you use these products, it behooves you to study their capabilities closely, compare them with the requirements of dimensional aggregation, test their application, and identify relevant implications. In fact, this is advised for users of *any* tool in Chapter 7.

MATERIALIZED VIEWS AND MATERIALIZED QUERY TABLES

The tool-agnostic principles and techniques in this book are periodically supplemented with a look at the impact of Oracle's *materialized views* and IBM DB/2's *materialized query tables*.

TOPIC	CHAPTER	DESCRIPTION
Aggregate Design	Chapter 3	Schema designers targeting these technologies should model the hierarchies implicit within a dimension table, and the relationships among their attributes. This information should be included in design documentation, along with defining queries for each aggregate table.
Aggregate Use	Chapter 4	Chapter 4 describes the use of materialized views and materialized query tables to provide query rewrite capabilities. It also shows how these technologies are used to implement aggregate fact tables, virtual aggregate dimension tables, and pre-joined aggregates.
Aggregate Refresh	Chapter 5	Materialized views and materialized query tables do not eliminate the need to manage the refresh of aggregates. It is also necessary to coordinate the refresh mechanism with the query rewrite mechanism.
Aggregate Construction	Chapter 6	It is not necessary to build an ETL process to load a materialized view or materialized query table. Once their refresh is configured, the database will take care of this job. But some adjustments to the base schema's ETL process may improve the overall performance of the aggregates.

Purpose of Each Chapter

This book is organized into chapters that address the major activities involved in the implementation of star schema aggregates. After establishing some fundamentals, chapters are dedicated to aggregate selection, design, usage, and construction. The remaining chapters address the organization of these activities into project plans, explore advanced design considerations, and address other impacts on the data warehouse.

Chapter 1: Fundamentals of Aggregates

This chapter establishes a foundation on which the rest of the book will build. It introduces the *star schema*, *aggregate tables*, and the *aggregate navigator*. Even if you are already familiar with these concepts, you should read Chapter 1. It establishes guiding principles for the development of *invisible aggregates*, which have zero impact on production applications. These principles will shape the best practices developed through the rest of the book. This chapter also introduces several forms of summarization that are not invisible to applications but may provide useful performance benefits.

Chapter 2: Choosing Aggregates

Chapter 2 takes on the difficult process of determining which aggregates should be built. You will learn how to identify and describe potential aggregates and determine the appropriate combination for implementation. This will require balancing the performance of potential aggregates with their potential usage and available resources. A variety of techniques will prove useful in identifying high-value aggregate tables.

Chapter 3: Designing Aggregates

The design of aggregate tables requires the same rigor as that of the base schema. Chapter 3 lays out a detailed set of principles for the design of dimensional aggregates. Best practices are identified and explained in detail, and a concrete set of deliverables is developed for the design process. Common pitfalls that can disrupt accuracy or ease of use are fully explored.

Chapter 4: Using Aggregates

In the most successful implementations, aggregate tables are invisible to users and applications. The job of the aggregate navigator is to redirect all queries to the best performing summaries. Chapter 4 develops a set of requirements for the aggregate navigator and uses them to evaluate three common styles of solutions. It explores two specific technologies in detail—Oracle’s materialized views and IBM DB/2’s materialized query tables—and provides practical advice for working without an aggregate navigator.

Chapter 5: ETL Part 1: Incorporating Aggregates

This book dedicates two chapters to the process of building aggregate tables. Chapter 5 describes how the base schema is loaded and how aggregates are integrated into that process. You will learn when it makes sense to design an incremental load for aggregate tables, and when you are better off dropping and rebuilding them each time the base schema is updated. For data warehouses loaded during batch windows, this chapter outlines several benefits of loading aggregates after the base schema. The ETL process will be required to interact with the aggregate navigator, or to take the entire data warehouse offline during the load. Data warehouses loaded in real-time require a different strategy for the maintenance of aggregates; specific techniques are discussed to minimize the impact of aggregates on this process.

Database features such as materialized views or materialized query tables may automate the construction process but are subject to the same requirements. As Chapter 5 shows, they must be configured to remain synchronized with the base schema, and designers must still choose between drop-and-rebuild and incremental load.

Chapter 6: ETL Part 2: Loading Aggregates

The second of two chapters on ETL, Chapter 6 describes the specific tasks required to load aggregate tables. Best practices are provided for identifying changed data in the base schema, constructing aggregate dimensions and their surrogate keys, and building aggregate fact tables. Pre-joined aggregates are also considered, along with complications that can arise from the presence of type 1 attributes.

The best practices in this chapter apply whether the load is developed using an ETL tool, or hand-coded. Database features such as materialized views or materialized query tables eliminate the need to design load routines, but may benefit from some adjustment to the schema design.

Chapter 7: Aggregates and Your Project

Aggregates should always be designed and implemented as part of a project. Chapter 7 provides a standard set of tasks and deliverables that can be used to add aggregates to existing schema, or to incorporate aggregates into the scope of a larger data warehouse development project. Major project phases are covered, including strategy, design, construction, testing, and deployment. The ongoing maintenance of aggregates is discussed, tying specific responsibilities to established data warehousing roles.

Chapter 8: Advanced Aggregate Design

This chapter outlines numerous advanced techniques for star schema design and fully analyzes the implications of each technique on aggregation. Design topics include:

- The periodic snapshot
- The accumulating snapshot
- Two kinds of factless fact tables
- Three kinds of bridge tables
- The transaction dimension
- Families of core and custom schemas

Chapter 8 also looks at how the techniques in this book can be adapted for snowflake schemas and third normal form designs.

Chapter 9: Related Topics

This final chapter collects several remaining topics that are influenced by aggregates:

- *The archive process* must be extended to involve aggregate tables. Some common misconceptions are discussed, and often-overlooked opportunities are highlighted.
- *Security requirements* may call for special care in implementing aggregates, which may also prove part of the solution.
- *Derived tables* are summarizations of base schema data that are not invisible. They include merged fact tables, sliced fact tables, and pivoted fact tables. Standard invisible aggregates may further summarize derived tables.
- *Deploying summary data before detail* can present new challenges, particularly if unanticipated. This chapter concludes by providing alternative techniques to deal with this unusual problem.

Glossary

Important terms used throughout this book are collected and defined in the glossary. You may find it useful to refer to these definitions as you read this book, particularly if you choose to read the chapters out of sequence.

Mastering Data Warehouse Aggregates

Fundamentals of Aggregates

A decade ago, Ralph Kimball described aggregate tables as “the single most dramatic way to improve performance in a large data warehouse.” Writing in *DBMS Magazine* (“Aggregate Navigation with (Almost) No Metadata,” August 1996), Kimball continued:

Aggregates can have a very significant effect on performance, in some cases speeding queries by a factor of one hundred or even one thousand. No other means exist to harvest such spectacular gains.

This statement rings as true today as it did ten years ago. Since then, advances in hardware and software have dramatically improved the capacity and performance of the data warehouse. Aggregates *compound* the effect of these improvements, providing performance gains that fully harness capabilities of the underlying technologies.

And the pressure to improve data warehouse performance is as strong as ever. As the baseline performance of underlying technologies has improved, warehouse developers have responded by storing and analyzing larger and more granular volumes of data. At the same time, warehouse systems have been opened to larger numbers of users, internal and external, who have come to expect instantaneous access to information.

This book empowers *you* to address these pressures. Using aggregate tables, you can achieve an extraordinary improvement in the speed of *your* data warehouse. And you can do it today, without making expensive upgrades to hardware, converting to a new database platform, or investing in exotic and proprietary technologies.

Although aggregates can have a powerful impact on data warehouse performance, they can also be misused. If not managed carefully, they can cause confusion, impose inordinate maintenance requirements, consume massive amounts of storage, and even provide inaccurate results. By following the best practices developed in this book, you can avoid these outcomes and maximize the positive impact of aggregates.

The introduction of aggregate tables to the data warehouse will touch every aspect of the data warehouse lifecycle. A set of best practices governs their selection, design, construction, and usage. They will influence data warehouse planning, project scope, maintenance requirements, and even the archive process. Before exploring each of these topics, it is necessary to establish some fundamental principles and vocabulary.

This chapter establishes the foundation on which the rest of the book builds. It introduces the *star schema*, *aggregate tables*, and the *aggregate navigator*. Guiding principles are established for the development of *invisible aggregates*, which have zero impact on production applications—other than performance, of course. Last, this chapter explores several other forms of summarization that are not invisible to applications, but may also provide useful performance benefits.

Star Schema Basics

A star schema is a set of tables in a relational database that has been designed according to the principles of *dimensional modeling*. Ralph Kimball popularized this approach to data warehouse design in the 1990s. Through his work and writings, Kimball established standard terminology and best practices that are now used around the world to design and build data warehouse systems. With coauthor Margy Ross, he provides a detailed treatment of these principles in *The Data Warehouse Toolkit, Second Edition* (Wiley, 2002).

To follow the examples throughout this book, you must understand the fundamental principles of dimensional modeling. In particular, the reader must have a basic grasp of the following concepts:

- The differences between data warehouse systems and operational systems
- How facts and dimensions support the measurement of a business process
- The tables of a star schema (fact tables and dimension tables) and their purposes

- The purpose of surrogate keys in dimension tables
- The grain of a fact table
- The additivity of facts
- How a star schema is queried
- Drilling across multiple fact tables
- Conformed dimensions and the warehouse bus
- The basic architecture of a data warehouse, including ETL software and BI software

If you are familiar with these topics, you may wish to skip to the section “Invisible Aggregates,” later in this chapter.

For everyone else, this section will bring you up-to-speed. Although not a substitute for Kimball and Ross’s book, this overview provides the background needed to understand the examples throughout this book. I encourage *all* readers to read *The Toolkit* for more immersion in the principles of dimensional modeling, particularly anyone involved in the design of the dimensional data warehouse.

Data warehouse designers will also benefit from reading *Data Warehouse Design Solutions*, by Chris Adamson and Mike Venerable (Wiley, 1998). This book explores the application of these principles in the service of specific business objectives and covers standard business processes in a wide variety of industries.

Operational Systems and the Data Warehouse

Data warehouse systems and operational systems have fundamentally different purposes. An operational system supports the *execution* of business process, while the data warehouse supports the *evaluation* of the process. Their distinct purposes are reflected in contrasting usage profiles, which in turn suggest that different principles will guide their design. The principles of dimensional modeling are specifically adapted to the unique requirements of the warehouse system.

Operational Systems

An operational system directly supports the *execution* of business processes. By capturing detail about significant events or transactions, it constructs a record of the activity. A sales system, for example, captures information about orders, shipments, and returns; a human resources system captures information about the hiring and promotion of employees; an accounting system captures information about the management of the financial assets and liabilities of the business. Capturing the detail surrounding these activities is often so important that the operational system becomes a part of the process.