

# **Data-Centric Systems and Applications**

## **Series Editors**

Michael J. Carey, University of California, Irvine, CA, USA

Stefano Ceri, Politecnico di Milano, Milano, Italy

## **Editorial Board Members**

Anastasia Ailamaki, École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

Shivnath Babu, Duke University, Durham, NC, USA

Philip A. Bernstein, Microsoft Corporation, Redmond, WA, USA

Johann-Christoph Freytag, Humboldt Universität zu Berlin, Berlin, Germany

Alon Halevy, Facebook, Menlo Park, CA, USA

Jiawei Han, University of Illinois, Urbana, IL, USA

Donald Kossmann, Microsoft Research Laboratory, Redmond, WA, USA

Gerhard Weikum, Max-Planck-Institut für Informatik, Saarbrücken, Germany

Kyu-Young Whang, Korea Advanced Institute of Science & Technology, Daejeon, Korea (Republic of)

Jeffrey Xu Yu, Chinese University of Hong Kong, Shatin, Hong Kong

Intelligent data management is the backbone of all information processing and has hence been one of the core topics in computer science from its very start. This series is intended to offer an international platform for the timely publication of all topics relevant to the development of data-centric systems and applications. All books show a strong practical or application relevance as well as a thorough scientific basis. They are therefore of particular interest to both researchers and professionals wishing to acquire detailed knowledge about concepts of which they need to make intelligent use when designing advanced solutions for their own problems.

Special emphasis is laid upon:

- Scientifically solid and detailed explanations of practically relevant concepts and techniques
  - (what does it do)
- Detailed explanations of the practical relevance and importance of concepts and techniques
  - (why do we need it)
- Detailed explanation of gaps between theory and practice
  - (why it does not work)

According to this focus of the series, submissions of advanced textbooks or books for advanced professional use are encouraged; these should preferably be authored books or monographs, but coherently edited, multi-author books are also envisaged (e.g. for emerging topics). On the other hand, overly technical topics (like physical data access, data compression etc.), latest research results that still need validation through the research community, or mostly product-related information for practitioners (“how to use Oracle 9i efficiently”) are not encouraged.

Alejandro Vaisman • Esteban Zimányi

# Data Warehouse Systems

Design and Implementation

Second Edition



Springer

Alejandro Vaisman    
Instituto Tecnológico de Buenos Aires  
Buenos Aires, Argentina

Esteban Zimányi    
Université Libre de Bruxelles  
Brussels, Belgium

ISSN 2197-9723                    ISSN 2197-974X (electronic)  
Data-Centric Systems and Applications  
ISBN 978-3-662-65166-7        ISBN 978-3-662-65167-4 (eBook)  
<https://doi.org/10.1007/978-3-662-65167-4>

© Springer-Verlag GmbH Germany, part of Springer Nature 2014, 2022

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

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

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer-Verlag GmbH, DE part of Springer Nature.

The registered company address is: Heidelberger Platz 3, 14197 Berlin, Germany

*To Andrés and Manuel,  
who bring me joy and  
happiness day after day*

*A. V.*

*To Elena,  
the star that shed light upon my path,  
with all my love*

*E.Z.*

# Foreword to the Second Edition

Dear reader,

Assuming you are looking for a textbook on data warehousing and the analytical processing of data, I can assure you that you are certainly in the right spot. In fact, I could easily argue how panoramic and lucid the view from this spot is, and in the next few paragraphs, this is exactly what I am going to do.

Assembling a good book from the bits and pieces of writings, slides, and article commentaries that an author has in his folders, is no easy task. Even more, if the book is intended to serve as a textbook, it requires an extra dose of love and care for the students who are going to use it (and their instructors, too, in fact). The book you have at hand is the product of hard work and deep caring by our two esteemed colleagues, Alejandro Vaisman and Esteban Zimányi, who have invested a large amount of effort to produce a book that is (a) comprehensive, (b) up-to-date, (c) easy to follow, and, (d) useful and to-the-point. While the book is also addressing the researcher who, coming from a different background, wants to enter the area of data warehousing, as well as the newcomer to data processing, who might prefer to start the journey of working with data from the neat setup of data cubes, the book is perfectly suited as a textbook for advanced undergraduate and graduate courses in the area of data warehousing.

The book comprehensively covers all the fundamental modeling issues, and addresses also the practical aspects on querying and populating the warehouse. The usage of concrete examples, consistently revisited throughout the book, guide the student to understand the practical considerations, and a set of exercises help the instructor with the hands-on design of a course. For what it's worth, I have already used the first edition of the book for my graduate data warehouse course and will certainly switch to the new version in the years to come.

If you, dear reader, have already read the first edition of the book, you already know that the first part, covering the modeling fundamentals, and the second part, covering the practical usage of data warehousing are both

comprehensive and detailed. To the extent that the fundamentals have not changed (and are not really expected to change in the future), apart from a set of extensions spread throughout the first part of the book, the main improvements concern readability on the one hand, and the technological advances on the other. Specifically, the dedicated chapter 7 on practical data analysis with lots of examples over a specific example, as well as the new topics covering partitioning and parallel data processing in the physical management of the data warehouse provide an even more easy path to the novice reader into the areas of querying and managing the warehouse.

I would like, however, to take the opportunity and direct your attention to the really new features of this second edition, which are found in the last unit of the book, concerning advanced areas of data warehousing. This part goes beyond the traditional data warehousing modeling and implementation and is practically completely refreshed compared to the first edition of the book. The chapter on temporal and multiversion warehousing covers the problem of time encoding for evolving facts and the management of versions. The part on spatial warehouses has been significantly updated. There is a brand-new chapter on graph data processing, and its application to graph warehousing and graph OLAP. Last but extremely significant, the crown jewel of the book, a brand-new chapter on the management of Big Data and the usage of Hadoop, Spark and Kylin, as well as the coverage of distributed, in-memory, columnar, and Not-Only-SQL DBMS's in the context of analytical data processing. Recent advents like data processing in the cloud, polystores and data lakes are also covered in the chapter.

Based on all that, dear reader, I can only invite you to dive into the contents of the book, feeling certain that, once you have completed its reading (or maybe, targeted parts of it), you will join me in expressing our gratitude to Alejandro and Esteban, for providing such a comprehensive textbook for the field of data warehousing in the first place, and for keeping it up to date with the recent developments, in this, current, second edition.

Ioannina, Greece

Panos Vassiliadis

# Foreword to the First Edition

Having worked with data warehouses for almost 20 years, I was both honored and excited when two veteran authors in the field asked me to write a foreword for their new book and sent me a PDF file with the current draft. Already the size of the PDF file gave me a first impression of a very comprehensive book, an impression that was heavily reinforced by reading the Table of Contents. After reading the entire book, I think it is quite simply the most comprehensive textbook about data warehousing on the market.

The book is very well suited for one or more data warehouse courses, ranging from the most basic to the most advanced. It has all the features that are necessary to make a good textbook. First, a running case study, based on the Northwind database known from Microsoft's tools, is used to illustrate all aspects using many detailed figures and examples. Second, key terms and concepts are highlighted in the text for better reading and understanding. Third, review questions are provided at the end of each chapter so students can quickly check their understanding. Fourth, the many detailed exercises for each chapter put the presented knowledge into action, yielding deep learning and taking students through all the steps needed to develop a data warehouse. Finally, the book shows how to implement data warehouses using leading industrial and open-source tools, concretely Microsoft's suite of data warehouse tools, giving students the essential hands-on experience that enables them to put the knowledge into practice.

For the complete database novice, there is even an introductory chapter on standard database concepts and design, making the book self-contained even for this group. It is quite impressive to cover all this material, usually the topic of an entire textbook, without making it a dense read. Next, the book provides a good introduction to basic multidimensional concepts, later moving on to advanced concepts such as summarizability. A complete overview of the data warehouse and online analytical processing (OLAP) "architecture stack" is given. For the conceptual modeling of the data warehouse, a concise and intuitive graphical notation is used, a full specification of which is given in

an appendix, along with a methodology for the modeling and the translation to (logical-level) relational schemas.

Later, the book provides a lot of useful knowledge about designing and querying data warehouses, including a detailed, yet easy to read, description of the de facto standard OLAP query language: MultiDimensional eXpressions (MDX). I certainly learned a thing or two about MDX in a short time. The chapter on extract-transform-load (ETL) takes a refreshingly different approach by using a graphical notation based on the Business Process Modeling Notation (BPMN), thus treating the ETL flow at a higher and more understandable level. Unlike most other data warehouse books, this book also provides comprehensive coverage on analytics, including data mining and reporting, and on how to implement these using industrial tools. The book even has a chapter on methodology issues such as requirements capture and the data warehouse development process, again something not covered by most data warehouse textbooks.

However, the one thing that really sets this book apart from its peers is the coverage of advanced data warehouse topics, such as spatial databases and data warehouses, spatiotemporal or mobility databases and data warehouses, and semantic web data warehouses. The book also provides a useful overview of novel “big data” technologies like Hadoop and novel database and data warehouse architectures like in-memory database systems, column store systems, and right-time data warehouses. These advanced topics are a distinguishing feature not found in other textbooks.

Finally, the book concludes by pointing to a number of exciting directions for future research in data warehousing, making it an interesting read even for seasoned data warehouse researchers.

A famous quote by IBM veteran Bruce Lindsay states that “relational databases are the foundation of Western civilization.” Similarly, I would say that “data warehouses are the foundation of twenty-first-century enterprises.” And this book is in turn an excellent foundation for building those data warehouses, from the simplest to the most complex.

Happy reading!

Aalborg, Denmark

Torben Bach Pedersen

# Preface

Since the late 1970s, relational database technology has been adopted by most organizations to store their essential data. However, nowadays, the needs of these organizations are not the same as they used to be. On the one hand, increasing market dynamics and competitiveness led to the need to have the right information at the right time. Managers need to be properly informed in order to take appropriate decisions to keep up with business successfully. On the other hand, data held by organizations are usually scattered among different systems, each one devised for a particular kind of business activity. Further, these systems may also be distributed geographically in different branches of the organization.

Traditional database systems are not well suited for these new requirements, since they were devised to support day-to-day operations rather than for data analysis and decision making. As a consequence, new database technologies for these specific tasks emerged in the 1990s, namely, data warehousing and online analytical processing (OLAP), which involve architectures, algorithms, tools, and techniques for bringing together data from heterogeneous information sources into a single repository suited for analysis. In this repository, called a data warehouse, data are accumulated over a period of time for the purpose of analyzing their evolution and discovering strategic information such as trends, correlations, and the like. Data warehousing is a well-established and mature technology used by organizations to improve their operations and better achieve their objectives.

## Objective of the Book

This book is aimed at consolidating and transferring to the community the experience of many years of teaching and research in the field of databases and data warehouses conducted by the authors, individually as well as jointly. However, this is not a compilation of the authors' past publications. On the

contrary, the book aims at being a main textbook for undergraduate and graduate computer science courses on data warehousing and OLAP. As such, it is written in a pedagogical rather than research style to make the work of the instructor easier and to help the student understand the concepts being delivered. Researchers and practitioners who are interested in an introduction to the area of data warehousing will also find in the book a useful reference. In summary, we aim at providing in-depth coverage of the main topics in the field, yet keeping a simple and understandable style.

Throughout the book, we cover all the phases of the data warehousing process, from requirements specification to implementation. Regarding data warehouse design, we make a clear distinction between the three abstraction levels of the American National Standards Institute (ANSI) database architecture, that is, conceptual, logical, and physical, unlike the usual approaches, which do not distinguish clearly between the conceptual and logical levels. A strong emphasis is placed on querying using the de facto standard language MDX (MultiDimensional eXpressions) as well as the popular language DAX (Data Analysis eXpressions). Though there are many practical books covering these languages, academic books have largely ignored them. We also provide in-depth coverage of the extraction, transformation, and loading (ETL) processes. In addition, we study how key performance indicators (KPIs) and dashboards are built on top of data warehouses. An important topic that we also cover in this book is temporal and multiversion data warehouses, in which the evolution over time of the data and the schema of a data warehouse are taken into account. Although there are many textbooks on spatial databases, this is not the case with spatial data warehouses, which we study in this book, together with mobility data warehouses, which allow the analysis of data produced by objects that change their position in space and time, like cars or pedestrians. Data warehousing and OLAP on graph databases and on the semantic web are also studied. Finally, big data technologies led to the concept of big data warehouses, which are also covered in this book.

A key characteristic that distinguishes this book from other textbooks is that we illustrate how the concepts introduced can be implemented using existing tools. Specifically, throughout the book we develop a case study based on the well-known Northwind database using representative tools of different kinds. In particular, the chapter on logical design includes a complete description of how to define an OLAP cube in Microsoft SQL Analysis Services using both the multidimensional and the tabular models. Similarly, the chapter on physical design illustrates how to optimize SQL Server and Analysis Services applications. Further, in the chapter on ETL we give a complete example of a process that loads the Northwind data warehouse, implemented using Integration Services. We also use Analysis Services for defining KPIs, and use Reporting Services to show how dashboards can be implemented. To illustrate spatial and spatiotemporal concepts we use the open-source database PostgreSQL, its spatial extension PostGIS, and its mobility extension MobilityDB. In this way, the reader can replicate most of the examples and queries

presented in the book. Finally, in the chapter on graph data warehouses we use Neo4j.

We also include review questions and exercises for all the chapters in order to help the reader verify that the concepts have been well understood. Support material for the book is available online at <http://cs.ulb.ac.be/DWSDIbook2e/>. This includes electronic versions of the figures, slides for each chapter, solutions to the exercises, and other pedagogic material that can be used by instructors using this book as a course text.

This second edition of the book updates several chapters with new results and technologies that have appeared since the publication of the first edition. In Chaps. 5, 6, and 7, the tabular model and DAX have been included. Chapter 15 covers big data warehouse technologies, which have considerably evolved since the first edition. Further, we have added new chapters covering temporal, multiversion, and graph data warehouses. Also, all application examples that make use of software tools have been updated to the latest versions of them. In addition to this new material, all chapters of the first edition have been revised and updated with the feedback obtained through seven years of teaching at undergraduate and graduate levels, and to professional teams in different industries.

## Organization of the Book and Teaching Paths

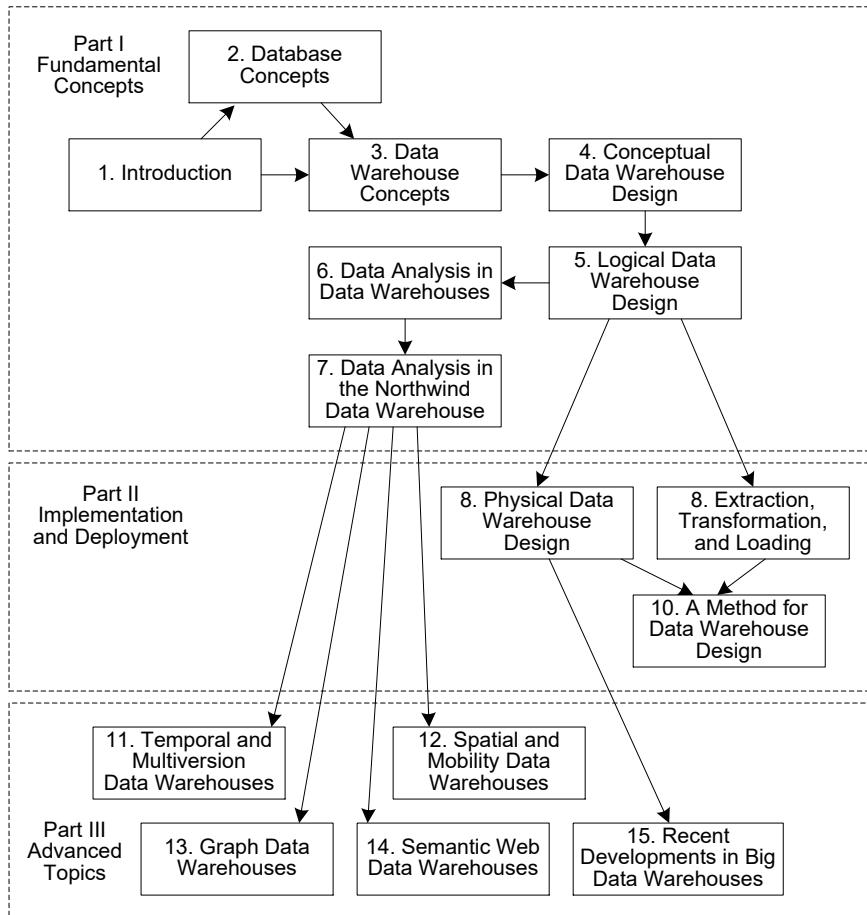
Part I of the book starts with Chap. 1, giving a historical overview of data warehousing and OLAP. Chapter 2 introduces the main concepts of relational databases needed in the remainder of the book. We also introduce the case study that we will use throughout the book, based on the well-known Northwind database. Data warehouses and the multidimensional model are introduced in Chap. 3, as well as the suite of tools provided by SQL Server. Chapter 4 deals with conceptual data warehouse design, while Chap. 5 is devoted to logical data warehouse design. Part I closes with Chaps. 6 and 7, which study SQL/OLAP, the extension of SQL with OLAP features, as well as MDX and DAX.

Part II covers data warehouse implementation issues. This part starts with Chap. 8, which tackles classical physical data warehouse design, focusing on indexing, view materialization, and database partitioning. Chapter 9 studies conceptual modeling and implementation of ETL processes. Finally, Chap. 10 provides a comprehensive method for data warehouse design.

Part III covers advanced data warehouse topics. This part starts with Chap. 11, which studies temporal and multiversion data warehouses, for both *data* and *schema* evolution of the data warehouse. Then, in Chap. 12, we study spatial data warehouses and their exploitation, denoted spatial OLAP (SOLAP), illustrating the problem with a spatial extension of the Northwind data warehouse denoted GeoNorthwind. We query this data warehouse

using PostGIS, PostgreSQL's spatial extension. The chapter also covers mobility data warehousing, using MobilityDB, a spatiotemporal extension of PostgreSQL. Chapters 13 and 14 address OLAP analysis over graph data represented, respectively, natively using property graphs in Neo4j and using RDF triples as advocated by the semantic web. Chapter 15 studies how novel techniques and technologies for distributed data storage and processing can be applied to the field of data warehousing. Appendix A summarizes the notations used in this book.

The figure below illustrates the overall structure of the book and the inter-dependencies between the chapters described above. Readers may refer to this figure to tailor their use of this book to their own particular interests. The dependency graph in the figure suggests many of the possible combinations that can be devised to offer advanced graduate courses on data warehousing.



Relationships between the chapters of this book

# Acknowledgments

We would like to thank Innoviris, the Brussels Institute for Research and Innovation, which funded Alejandro Vaisman's work through the OSCB project; without its financial support, the first edition of this book would never have been possible. As mentioned above, some content of this book finds its roots in a previous book written by one of the authors in collaboration with Elzbieta Malinowski. We would like to thank her for all the work we did together in making the previous book a reality. This gave us the impetus to start this new book.

Parts of the material included in this book have been previously presented in conferences or published in journals. At these conferences, we had the opportunity to discuss with research colleagues from all around the world, and we exchanged viewpoints about the subject with them. The anonymous reviewers of these conferences and journals provided us with insightful comments and suggestions that contributed significantly to improve the work presented in this book. We would like to thank Zineb El Akkaoui, with whom we have explored the use of BPMN for ETL processes, and Judith Awiti, who continued this work. A very special thanks to Waqas Ahmed, a doctoral student of our laboratory, with whom we explored the issue of temporal and multiversion data warehouses. Waqas also suggested to include tabular modeling and DAX in the second edition of the book, and without his invaluable help, all the material related to the tabular model and DAX would have not been possible. A special thanks to Mahmoud Sakr, Arthur Lesuisse, Mohammed Bakli, and Maxime Schoemans, who worked with one of the authors in the development of MobilityDB, a spatiotemporal extension of PostgreSQL and PostGIS that was used for mobility data warehouses. This work follows that of Benoit Foé, Julien Lusiela, and Xianling Li, who explored this topic in the context of their master's thesis. Arthur Lesuisse also provided invaluable help in setting up all the computer infrastructure we needed, especially for spatializing the Northwind database. He also contributed in enhancing some of the figures of this book. Thanks also to Leticia Gómez from the Buenos Aires Technological Institute for her help on the im-

plementation of graph data warehouses and for her advice on the topic of big data technologies. Bart Kuijpers, from Hasselt University, also worked with us during our research on graph data warehousing and OLAP. We also want to thank Lorena Etcheverry, who contributed with comments, exercises, and solutions in Chap. 14.

Special thanks go to Panos Vassiliadis, professor at the University of Ioannina in Greece, who kindly agreed to write the foreword for this second edition. Finally, we would like to warmly thank Ralf Gerstner of Springer for his continued interest in this book. The enthusiastic welcome given to our book proposal for the first edition and the continuous encouragements to write the second edition gave us enormous impetus to pursue our project to its end.

Alejandro Vaisman  
Buenos Aires, Argentina,  
February 2022

Esteban Zimányi  
Brussels, Belgium

# About the Authors

**Alejandro Vaisman** is a professor at the Instituto Tecnológico de Buenos Aires, where he also chairs the graduate program in data science. He has been a professor and chair of the master's program in data mining at the University of Buenos Aires (UBA) and professor at Universidad de la República in Uruguay. He received a BE degree in civil engineering, and a BCS degree and a doctorate in computer science from the UBA, under the supervision of Prof. Alberto Mendelzon, from the University of Toronto (UoT). He has been a postdoctoral fellow at UoT, and visiting researcher at UoT, Universidad Politécnica de Madrid, Universidad de Chile, University of Hasselt, and Université Libre de Bruxelles (ULB). His research interests are in the field of databases, business intelligence, and geographic information systems. He has authored and coauthored many scientific papers published at major conferences and in major journals.

**Esteban Zimányi** is a professor and a director of the Department of Computer and Decision Engineering (CoDE) of Université Libre de Bruxelles (ULB). He started his studies at the Universidad Autónoma de Centro América, Costa Rica, and received a BCS degree and a doctorate in computer science from ULB. His current research interests include spatiotemporal and mobility databases, data warehouses and business intelligence, geographic information systems, as well as semantic web. He has coauthored and coedited eight books and published many papers on these topics. He was editor-in-chief of the *Journal on Data Semantics* (JoDS) published by Springer from 2012 to 2020. He coordinated the Erasmus Mundus master's and doctorate programmes "Information Technologies for Business Intelligence" (IT4BI) and "Big Data Management and Analytics" (BDMA) as well as the Marie Skłodowska-Curie doctorate programme "Data Engineering for Data Science" (DEDS).

# Contents

## Part I Fundamental Concepts

|          |  |    |
|----------|--|----|
| <b>1</b> | <b>Introduction</b>                    | 3  |
| 1.1      | An Overview of Data Warehousing        | 4  |
| 1.2      | Emerging Data Warehousing Technologies | 7  |
| 1.3      | Review Questions                       | 10 |
| <b>2</b> | <b>Database Concepts</b>               | 11 |
| 2.1      | Database Design                        | 11 |
| 2.2      | The Northwind Case Study               | 13 |
| 2.3      | Conceptual Database Design             | 13 |
| 2.4      | Logical Database Design                | 18 |
| 2.4.1    | The Relational Model                   | 18 |
| 2.4.2    | Normalization                          | 24 |
| 2.4.3    | Relational Query Languages             | 26 |
| 2.5      | Physical Database Design               | 36 |
| 2.6      | Summary                                | 39 |
| 2.7      | Bibliographic Notes                    | 40 |
| 2.8      | Review Questions                       | 40 |
| 2.9      | Exercises                              | 41 |
| <b>3</b> | <b>Data Warehouse Concepts</b>         | 45 |
| 3.1      | Multidimensional Model                 | 45 |
| 3.1.1    | Hierarchies                            | 47 |
| 3.1.2    | Measures                               | 49 |
| 3.2      | OLAP Operations                        | 51 |
| 3.3      | Data Warehouses                        | 63 |
| 3.4      | Data Warehouse Architecture            | 67 |
| 3.4.1    | Back-End Tier                          | 67 |
| 3.4.2    | Data Warehouse Tier                    | 68 |
| 3.4.3    | OLAP Tier                              | 69 |

|          |   |            |
|----------|---|------------|
| 3.4.4    | Front-End Tier . . . . .  | 70         |
| 3.4.5    | Variations of the Architecture . . . . .                        | 70         |
| 3.5      | Overview of Microsoft SQL Server BI Tools . . . . .             | 71         |
| 3.6      | Summary . . . . .   | 72         |
| 3.7      | Bibliographic Notes . . . . .                                   | 72         |
| 3.8      | Review Questions . . . . .                                      | 73         |
| 3.9      | Exercises . . . . .   | 73         |
| <b>4</b> | <b>Conceptual Data Warehouse Design . . . . .</b>               | <b>75</b>  |
| 4.1      | Conceptual Modeling of Data Warehouses . . . . .                | 75         |
| 4.2      | Hierarchies . . . . .   | 79         |
| 4.2.1    | Balanced Hierarchies . . . . .                                  | 79         |
| 4.2.2    | Unbalanced Hierarchies . . . . .                                | 80         |
| 4.2.3    | Generalized Hierarchies . . . . .                               | 81         |
| 4.2.4    | Alternative Hierarchies . . . . .                               | 83         |
| 4.2.5    | Parallel Hierarchies . . . . .                                  | 84         |
| 4.2.6    | Nonstrict Hierarchies . . . . .                                 | 87         |
| 4.3      | Advanced Modeling Aspects . . . . .                             | 90         |
| 4.3.1    | Facts with Multiple Granularities . . . . .                     | 91         |
| 4.3.2    | Many-to-Many Dimensions . . . . .                               | 91         |
| 4.3.3    | Links between Facts . . . . .                                   | 95         |
| 4.4      | Querying the Northwind Cube Using the OLAP Operations . . . . . | 96         |
| 4.5      | Summary . . . . .   | 99         |
| 4.6      | Bibliographic Notes . . . . .                                   | 100        |
| 4.7      | Review Questions . . . . .                                      | 101        |
| 4.8      | Exercises . . . . .   | 102        |
| <b>5</b> | <b>Logical Data Warehouse Design . . . . .</b>                  | <b>105</b> |
| 5.1      | Logical Modeling of Data Warehouses . . . . .                   | 105        |
| 5.2      | Relational Data Warehouse Design . . . . .                      | 106        |
| 5.3      | Relational Representation of Data Warehouses . . . . .          | 109        |
| 5.4      | Time Dimension . . . . .  | 112        |
| 5.5      | Logical Representation of Hierarchies . . . . .                 | 113        |
| 5.5.1    | Balanced Hierarchies . . . . .                                  | 113        |
| 5.5.2    | Unbalanced Hierarchies . . . . .                                | 114        |
| 5.5.3    | Generalized Hierarchies . . . . .                               | 115        |
| 5.5.4    | Alternative Hierarchies . . . . .                               | 117        |
| 5.5.5    | Parallel Hierarchies . . . . .                                  | 118        |
| 5.5.6    | Nonstrict Hierarchies . . . . .                                 | 119        |
| 5.6      | Advanced Modeling Aspects . . . . .                             | 120        |
| 5.6.1    | Facts with Multiple Granularities . . . . .                     | 120        |
| 5.6.2    | Many-to-Many Dimensions . . . . .                               | 121        |
| 5.6.3    | Links between Facts . . . . .                                   | 122        |
| 5.7      | Slowly Changing Dimensions . . . . .                            | 124        |
| 5.8      | Performing OLAP Queries with SQL . . . . .                      | 130        |

|          |  |            |
|----------|--|------------|
| 5.9      | Defining the Northwind Data Warehouse in Analysis Services | 135        |
| 5.9.1    | Multidimensional Model                                     | 135        |
| 5.9.2    | Tabular Model  | 147        |
| 5.10     | Summary  | 155        |
| 5.11     | Bibliographic Notes  | 156        |
| 5.12     | Review Questions   | 156        |
| 5.13     | Exercises  | 157        |
| <b>6</b> | <b>Data Analysis in Data Warehouses</b>                    | <b>159</b> |
| 6.1      | Introduction to MDX  | 160        |
| 6.1.1    | Tuples and Sets  | 160        |
| 6.1.2    | Basic Queries  | 162        |
| 6.1.3    | Slicing  | 163        |
| 6.1.4    | Navigation   | 164        |
| 6.1.5    | Cross Join   | 166        |
| 6.1.6    | Subqueries   | 167        |
| 6.1.7    | Calculated Members and Named Sets                          | 168        |
| 6.1.8    | Relative Navigation  | 170        |
| 6.1.9    | Time-Related Calculations                                  | 171        |
| 6.1.10   | Filtering  | 174        |
| 6.1.11   | Sorting  | 175        |
| 6.1.12   | Top and Bottom Analysis                                    | 176        |
| 6.1.13   | Aggregation Functions                                      | 177        |
| 6.2      | Introduction to DAX  | 179        |
| 6.2.1    | Expressions  | 179        |
| 6.2.2    | Evaluation Context   | 181        |
| 6.2.3    | Queries  | 183        |
| 6.2.4    | Filtering  | 184        |
| 6.2.5    | Hierarchy Handling   | 188        |
| 6.2.6    | Time-Related Calculations                                  | 189        |
| 6.2.7    | Top and Bottom Analysis                                    | 192        |
| 6.2.8    | Table Operations   | 194        |
| 6.3      | Key Performance Indicators                                 | 196        |
| 6.3.1    | Classification of Key Performance Indicators               | 197        |
| 6.3.2    | Defining Key Performance Indicators                        | 198        |
| 6.4      | Dashboards   | 200        |
| 6.4.1    | Types of Dashboards  | 201        |
| 6.4.2    | Guidelines for Dashboard Design                            | 202        |
| 6.5      | Summary  | 203        |
| 6.6      | Bibliographic Notes  | 203        |
| 6.7      | Review Questions   | 204        |

|          |  |     |
|----------|--|-----|
| <b>7</b> | <b>Data Analysis in the Northwind Data Warehouse . . . . .</b> | 205 |
| 7.1      | Querying the Multidimensional Model in MDX . . . . .           | 205 |
| 7.2      | Querying the Tabular Model in DAX . . . . .                    | 211 |
| 7.3      | Querying the Relational Data Warehouse in SQL . . . . .        | 217 |
| 7.4      | Comparison of MDX, DAX, and SQL . . . . .                      | 225 |
| 7.5      | KPIs for the Northwind Case Study . . . . .                    | 229 |
| 7.5.1    | KPIs in Analysis Services Multidimensional . . . . .           | 229 |
| 7.5.2    | KPIs in Analysis Services Tabular . . . . .                    | 232 |
| 7.6      | Dashboards for the Northwind Case Study . . . . .              | 234 |
| 7.6.1    | Dashboards in Reporting Services . . . . .                     | 235 |
| 7.6.2    | Dashboards in Power BI . . . . .                               | 239 |
| 7.7      | Summary . . . . .  | 241 |
| 7.8      | Review Questions . . . . .                                     | 241 |
| 7.9      | Exercises . . . . .  | 242 |

## Part II Implementation and Deployment

|          |   |     |
|----------|---|-----|
| <b>8</b> | <b>Physical Data Warehouse Design . . . . .</b>               | 245 |
| 8.1      | Physical Modeling of Data Warehouses . . . . .                | 246 |
| 8.2      | Materialized Views . . . . .                                  | 247 |
| 8.2.1    | Algorithms Using Full Information . . . . .                   | 249 |
| 8.2.2    | Algorithms Using Partial Information . . . . .                | 251 |
| 8.3      | Data Cube Maintenance . . . . .                               | 252 |
| 8.4      | Computation of a Data Cube . . . . .                          | 258 |
| 8.4.1    | PipeSort Algorithm . . . . .                                  | 259 |
| 8.4.2    | Cube Size Estimation . . . . .                                | 262 |
| 8.4.3    | Partial Computation of a Data Cube . . . . .                  | 263 |
| 8.5      | Indexes for Data Warehouses . . . . .                         | 267 |
| 8.5.1    | Bitmap Indexes . . . . .                                      | 268 |
| 8.5.2    | Bitmap Compression . . . . .                                  | 271 |
| 8.5.3    | Join Indexes . . . . .  | 272 |
| 8.6      | Evaluation of Star Queries . . . . .                          | 273 |
| 8.7      | Partitioning . . . . .  | 275 |
| 8.8      | Parallel Processing . . . . .                                 | 277 |
| 8.9      | Physical Design in SQL Server and Analysis Services . . . . . | 280 |
| 8.9.1    | Indexed Views . . . . .                                       | 280 |
| 8.9.2    | Partition-Aligned Indexed Views . . . . .                     | 282 |
| 8.9.3    | Column-Store Indexes . . . . .                                | 283 |
| 8.9.4    | Partitions in Analysis Services . . . . .                     | 284 |
| 8.10     | Query Performance in Analysis Services . . . . .              | 286 |
| 8.11     | Summary . . . . .   | 289 |
| 8.12     | Bibliographic Notes . . . . .                                 | 290 |
| 8.13     | Review Questions . . . . .                                    | 290 |
| 8.14     | Exercises . . . . .   | 291 |

|  |     |
|--|-----|
| <b>9 Extraction, Transformation, and Loading . . . . .</b>       | 297 |
| 9.1 Business Process Modeling Notation . . . . .                 | 298 |
| 9.2 Conceptual ETL Design Using BPMN . . . . .                   | 303 |
| 9.3 Conceptual Design of the Northwind ETL Process . . . . .     | 306 |
| 9.4 SQL Server Integration Services . . . . .                    | 318 |
| 9.5 The Northwind ETL Process in Integration Services . . . . .  | 320 |
| 9.6 Implementing ETL Processes in SQL . . . . .                  | 326 |
| 9.7 Summary . . . . .  | 332 |
| 9.8 Bibliographic Notes . . . . .                                | 332 |
| 9.9 Review Questions . . . . .                                   | 333 |
| 9.10 Exercises . . . . .   | 334 |
| <b>10 A Method for Data Warehouse Design . . . . .</b>           | 335 |
| 10.1 Approaches to Data Warehouse Design . . . . .               | 335 |
| 10.2 General Overview of the Method . . . . .                    | 337 |
| 10.3 Requirements Specification . . . . .                        | 338 |
| 10.3.1 Business-Driven Requirements Specification . . . . .      | 339 |
| 10.3.2 Data-driven Requirements Specification . . . . .          | 345 |
| 10.3.3 Business/Data-driven Requirements Specification . . . . . | 349 |
| 10.4 Conceptual Design . . . . .                                 | 350 |
| 10.4.1 Business-Driven Conceptual Design . . . . .               | 351 |
| 10.4.2 Data-driven Conceptual Design . . . . .                   | 354 |
| 10.4.3 Business/Data-driven Conceptual Design . . . . .          | 356 |
| 10.5 Logical Design . . . . .                                    | 357 |
| 10.5.1 Logical Schemas . . . . .                                 | 358 |
| 10.5.2 ETL Processes . . . . .                                   | 359 |
| 10.6 Physical Design . . . . .                                   | 359 |
| 10.7 Characterization of the Various Approaches . . . . .        | 360 |
| 10.7.1 Business-Driven Approach . . . . .                        | 360 |
| 10.7.2 Data-driven Approach . . . . .                            | 361 |
| 10.7.3 Business/Data-driven Approach . . . . .                   | 362 |
| 10.8 Summary . . . . .   | 363 |
| 10.9 Bibliographic Notes . . . . .                               | 363 |
| 10.10 Review Questions . . . . .                                 | 365 |
| 10.11 Exercises . . . . .  | 366 |

## Part III Advanced Topics

|  |     |
|--|-----|
| <b>11 Temporal and Multiversion Data Warehouses . . . . .</b>    | 373 |
| 11.1 Manipulating Temporal Information in SQL . . . . .          | 374 |
| 11.2 Conceptual Design of Temporal Data Warehouses . . . . .     | 383 |
| 11.2.1 Time Data Types . . . . .                                 | 383 |
| 11.2.2 Synchronization Relationships . . . . .                   | 384 |
| 11.2.3 A Conceptual Model for Temporal Data Warehouses . . . . . | 386 |
| 11.2.4 Temporal Hierarchies . . . . .                            | 389 |

|           |  |            |
|-----------|--|------------|
| 11.2.5    | Temporal Facts . . . . .   | 391        |
| 11.3      | Logical Design of Temporal Data Warehouses . . . . .                 | 392        |
| 11.4      | Implementation Considerations . . . . .                              | 395        |
| 11.4.1    | Period Encoding . . . . .  | 395        |
| 11.4.2    | Tables for Temporal Roll-Up . . . . .                                | 395        |
| 11.4.3    | Integrity Constraints . . . . .                                      | 396        |
| 11.4.4    | Measure Aggregation . . . . .  | 399        |
| 11.4.5    | Temporal Measures . . . . .  | 403        |
| 11.5      | Querying the Temporal Northwind Data Warehouse in SQL .              | 404        |
| 11.6      | Temporal Data Warehouses versus Slowly Changing Dimensions . . . . . | 412        |
| 11.7      | Conceptual Design of Multiversion Data Warehouses . . . . .          | 416        |
| 11.8      | Logical Design of Multiversion Data Warehouses . . . . .             | 422        |
| 11.9      | Querying the Multiversion Northwind Data Warehouse in SQL . . . . .  | 427        |
| 11.10     | Summary . . . . .  | 428        |
| 11.11     | Bibliographic Notes . . . . .  | 429        |
| 11.12     | Review Questions . . . . .   | 430        |
| 11.13     | Exercises . . . . .  | 431        |
| <b>12</b> | <b>Spatial and Mobility Data Warehouses . . . . .</b>                | <b>437</b> |
| 12.1      | Conceptual Design of Spatial Data Warehouses . . . . .               | 438        |
| 12.1.1    | Spatial Data Types . . . . .   | 438        |
| 12.1.2    | Topological relationships . . . . .                                  | 440        |
| 12.1.3    | Continuous Fields . . . . .  | 441        |
| 12.1.4    | A Conceptual Model of Spatial Data Warehouses . .                    | 441        |
| 12.2      | Implementation Considerations for Spatial Data . . . . .             | 445        |
| 12.2.1    | Spatial Reference Systems . . . . .                                  | 445        |
| 12.2.2    | Vector Model . . . . .   | 447        |
| 12.2.3    | Raster Model . . . . .   | 449        |
| 12.3      | Logical Design of Spatial Data Warehouses . . . . .                  | 451        |
| 12.4      | Topological Constraints . . . . .                                    | 454        |
| 12.5      | Querying the GeoNorthwind Data Warehouse in SQL . .                  | 456        |
| 12.6      | Mobility Data Analysis . . . . .                                     | 460        |
| 12.7      | Temporal Types . . . . .   | 461        |
| 12.8      | Temporal Types in MobilityDB . . . . .                               | 466        |
| 12.9      | Mobility Data Warehouses . . . . .                                   | 470        |
| 12.10     | Querying the Northwind Mobility Data Warehouse in SQL .              | 474        |
| 12.11     | Summary . . . . .  | 480        |
| 12.12     | Bibliographic Notes . . . . .  | 480        |
| 12.13     | Review Questions . . . . .   | 481        |
| 12.14     | Exercises . . . . .  | 482        |

|  |     |
|--|-----|
| <b>13 Graph Data Warehouses . . . . .</b>                      | 487 |
| 13.1 Graph Data Models . . . . .                               | 488 |
| 13.2 Property Graph Database Systems . . . . .                 | 490 |
| 13.2.1 Neo4j . . . . .   | 492 |
| 13.2.2 Introduction to Cypher . . . . .                        | 493 |
| 13.2.3 Querying the Northwind Cube with Cypher . . . . .       | 501 |
| 13.3 OLAP on Hypergraphs . . . . .                             | 507 |
| 13.3.1 Operations on Hypergraphs . . . . .                     | 512 |
| 13.3.2 OLAP on Trajectory Graphs . . . . .                     | 516 |
| 13.4 Graph Processing Frameworks . . . . .                     | 520 |
| 13.4.1 Gremlin . . . . .                                       | 520 |
| 13.4.2 JanusGraph . . . . .                                    | 523 |
| 13.5 Bibliographic Notes . . . . .                             | 526 |
| 13.6 Review Questions . . . . .                                | 526 |
| 13.7 Exercises . . . . .                                       | 527 |
| <b>14 Semantic Web Data Warehouses . . . . .</b>               | 531 |
| 14.1 Semantic Web . . . . .                                    | 532 |
| 14.1.1 Introduction to RDF and RDFS . . . . .                  | 532 |
| 14.1.2 RDF Serializations . . . . .                            | 533 |
| 14.1.3 RDF Representation of Relational Data . . . . .         | 535 |
| 14.2 Introduction to SPARQL . . . . .                          | 539 |
| 14.2.1 SPARQL Basics . . . . .                                 | 540 |
| 14.2.2 SPARQL Semantics . . . . .                              | 543 |
| 14.3 RDF Representation of Multidimensional Data . . . . .     | 544 |
| 14.4 Representation of the Northwind Cube in QB4OLAP . . . . . | 547 |
| 14.5 Querying the Northwind Cube in SPARQL . . . . .           | 549 |
| 14.6 Summary . . . . .   | 557 |
| 14.7 Bibliographic Notes . . . . .                             | 557 |
| 14.8 Review Questions . . . . .                                | 558 |
| 14.9 Exercises . . . . .                                       | 559 |
| <b>15 Recent Developments in Big Data Warehouses . . . . .</b> | 561 |
| 15.1 Data Warehousing in the Age of Big Data . . . . .         | 562 |
| 15.2 Distributed Processing Frameworks . . . . .               | 563 |
| 15.2.1 Hadoop . . . . .  | 565 |
| 15.2.2 Hive . . . . .  | 567 |
| 15.2.3 Spark . . . . .   | 569 |
| 15.2.4 Comparison of Hadoop and Spark . . . . .                | 576 |
| 15.2.5 Kylin . . . . .   | 577 |
| 15.3 Distributed Database Systems . . . . .                    | 579 |
| 15.3.1 MySQL Cluster . . . . .                                 | 582 |
| 15.3.2 Citus . . . . .   | 585 |
| 15.4 In-Memory Database Systems . . . . .                      | 587 |
| 15.4.1 Oracle TimesTen . . . . .                               | 590 |

|                   |  |     |
|-------------------|--|-----|
| 15.4.2            | Redis  | 591 |
| 15.5              | Column-Store Database Systems                  | 592 |
| 15.5.1            | Vertica  | 595 |
| 15.5.2            | MonetDB  | 597 |
| 15.5.3            | Citus Columnar                                 | 598 |
| 15.6              | NoSQL Database Systems                         | 599 |
| 15.6.1            | HBase  | 600 |
| 15.6.2            | Cassandra                                      | 602 |
| 15.7              | NewSQL Database Systems                        | 606 |
| 15.7.1            | Cloud Spanner                                  | 607 |
| 15.7.2            | SAP HANA                                       | 607 |
| 15.7.3            | VoltDB   | 609 |
| 15.8              | Array Database Systems                         | 610 |
| 15.8.1            | Rasdaman                                       | 612 |
| 15.8.2            | SciDB  | 614 |
| 15.9              | Hybrid Transactional and Analytical Processing | 616 |
| 15.9.1            | SingleStore                                    | 617 |
| 15.9.2            | LeanXcale                                      | 618 |
| 15.10             | Polystores                                     | 619 |
| 15.10.1           | CloudMdsQL                                     | 620 |
| 15.10.2           | BigDAWG  | 621 |
| 15.11             | Cloud Data Warehouses                          | 622 |
| 15.12             | Data Lakes and Data Lakehouses                 | 624 |
| 15.13             | Future Perspectives                            | 628 |
| 15.14             | Summary  | 629 |
| 15.15             | Bibliographic Notes                            | 629 |
| 15.16             | Review Questions                               | 630 |
| <b>A</b>          | <b>Graphical Notation</b>                      | 633 |
| A.1               | Entity-Relationship Model                      | 633 |
| A.2               | Relational Model                               | 635 |
| A.3               | MultiDim Model for Data Warehouses             | 635 |
| A.4               | MultiDim Model for Spatial Data Warehouses     | 639 |
| A.5               | MultiDim Model for Temporal Data Warehouses    | 641 |
| A.6               | BPMN Notation for ETL                          | 643 |
| <b>References</b> |  | 647 |
| <b>Glossary</b>   |  | 667 |
| <b>Index</b>      |  | 685 |