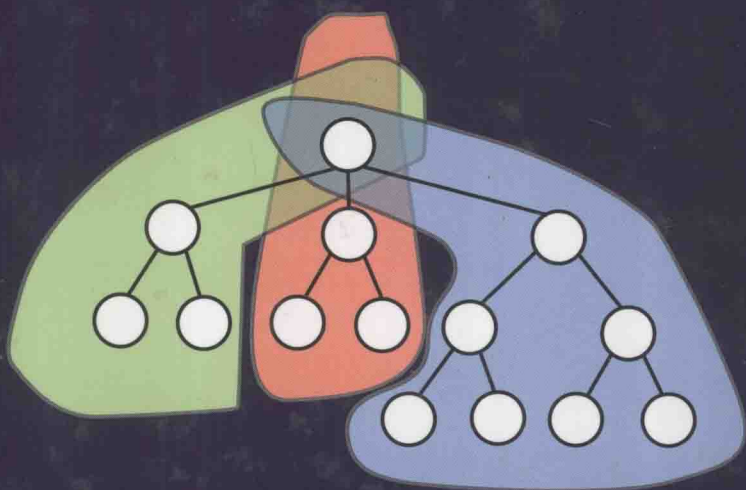# High-Performance Parallel Database Processing and Grid Databases

DAVID TANIAR • CLEMENT H. C. LEUNG

WENNY RAHAYU • SUSHANT GOEL

# High-Performance Parallel Database Processing and Grid Databases

**David Taniar**
*Monash University, Australia*

**Clement H.C. Leung**
*Hong Kong Baptist University and Victoria University, Australia*

**Wenny Rahayu**
*La Trobe University, Australia*

**Sushant Goel**
*RMIT University, Australia*

**WILEY**

# High-Performance
# Parallel Database
# Processing
# and Grid Databases

**WILEY SERIES ON PARALLEL
AND DISTRIBUTED COMPUTING**

**Editor: Albert Y. Zomaya**

A complete list of titles in this series appears at the end of this volume.

# Preface

The sizes of databases have seen exponential growth in the past, and such growth is expected to accelerate in the future, with the steady drop in storage cost accompanied by a rapid increase in storage capacity. Many years ago, a terabyte database was considered to be large, but nowadays they are sometimes regarded as small, and the daily volumes of data being added to some databases are measured in terabytes. In the future, petabyte and exabyte databases will be common.

With such volumes of data, it is evident that the sequential processing paradigm will be unable to cope; for example, even assuming a data rate of 1 terabyte per second, reading through a petabyte database will take over 10 days. To effectively manage such volumes of data, it is necessary to allocate multiple resources to it, very often massively so. The processing of databases of such astronomical proportions requires an understanding of how high-performance systems and parallelism work. Besides the massive volume of data in the database to be processed, some data has been distributed across the globe in a Grid environment. These massive data centers are also a part of the emergence of Cloud computing, where data access has shifted from local machines to powerful servers hosting web applications and services, making data access across the Internet using standard web browsers pervasive. This adds another dimension to such systems.

Parallelism in databases has been around since the early 1980s, when many researchers in this area aspired to build large special-purpose database machines—databases employing dedicated specialized parallel hardware. Some projects were born, including Bubba, Gamma, etc. These came and went. However, commercial DBMS vendors quickly realized the importance of supporting high performance for large databases, and many of them have incorporated parallelism and grid features into their products. Their commitment to high-performance systems and parallelism, as well as grid configurations, shows the importance and inevitability of parallelism.

In addition, while traditional transactional data is still common, we see an increasing growth of new application domains, broadly categorized as data-intensive applications. These include data warehousing and online analytic processing (OLAP) applications, data mining, genome databases, and multiple media databases manipulating unstructured and semistructured data. Therefore, it is critical to understand the underlying principle of data parallelism, before specialized and new application domains can be properly addressed.

This book is written to provide a fundamental understanding of parallelism in data-intensive applications. It features not only the algorithms for database operations but also quantitative analytical models, so that performance can be analyzed and evaluated more effectively.

The present book brings into a single volume the latest techniques and principles of parallel and grid database processing. It provides a much-needed, self-contained advanced text for database courses at the postgraduate or final year undergraduate levels. In addition, for researchers with a particular interest in parallel databases and related areas, it will serve as an indispensable and up-to-date reference. Practitioners contemplating building high-performance databases or seeking to gain a good understanding of parallel database technology too will find this book valuable for the wealth of techniques and models it contains.

## STRUCTURE OF THE BOOK

This book is divided into five parts. Part I gives an introduction to the topic, including the rationale behind the need for high-performance database processing, as well as basic analytical models that will be used throughout the book.

Part II, consisting of three chapters, describes parallelism for basic query operations. These include parallel searching, parallel aggregate and sorting, and parallel join. These are the foundation of query processing, whereby complex queries can be decomposed into any of these atomic operations.

Part III, consisting of the next four chapters, focuses on more advanced query operations. This part covers groupby-join operations, parallel indexing, parallel object-oriented query processing, in particular, collection join, and query scheduling and optimization.

Just as the previous two parts deal with parallelism of read-only queries, the next part, Part IV, concentrates on transactions, also known as write queries. We use the grid environment to study transaction management. In grid transaction management, the focus is mainly on grid concurrency control, atomic commitment, durability, as well as replication.

Finally, Part V introduces other data-intensive applications, including data warehousing, OLAP, business intelligence, and parallel data mining.

## ACKNOWLEDGMENTS

We also thank Bruna Pomella, who proofread the entire manuscript, for commenting on ambiguous sentences and correcting grammatical mistakes.

Finally, we would like to express our sincere thanks to our respective universities, Monash University, Victoria University, Hong Kong Baptist University, La Trobe University, and RMIT, where the research presented in this book was conducted. We are grateful for the facilities and time that we received during the writing of this book. Without these, the book would not have been written in the first place.

*David Taniar*
*Clement H.C. Leung*
*Wenny Rahayu*
*Sushant Goel*

# Contents

**Part III     Advanced Parallel Query Processing**

## 6. Parallel GroupBy-Join                        141

# Part I

# Introduction