# Querying Graphs

# Synthesis Lectures on Data Management

Querying Graphs
Angela Bonifati, George Fletcher, Hannes Voigt, and Nikolay Yakovets
2018

Query Processing over Incomplete Databases
Yunjun Gao and Xiaoye Miao
2018

Natural Language Data Management and Interfaces
Yunyao Li and Davood Rafiei
2018

Human Interaction with Graphs: A Visual Querying Perspective
Sourav S. Bhowmick, Byron Choi, and Chengkai Li
2018

On Uncertain Graphs
Arijit Khan, Yuan Ye, and Lei Chen
2018

Answering Queries Using Views
Foto Afrati and Rada Chirkova
2017

Querying Graphs

Angela Bonifati, George Fletcher, Hannes Voigt, and Nikolay Yakovets

# Querying Graphs

Angela Bonifati
Université Claude Bernard Lyon 1

George Fletcher
Technische Universiteit Eindhoven

Hannes Voigt
Neo4j/Technische Universität Dresden[1]

Nikolay Yakovets
Technische Universiteit Eindhoven

*SYNTHESIS LECTURES ON DATA MANAGEMENT #51*

[1]Author is now at Neo4j. The book was mainly written while the author was still at Technische Universität Dresden.

## ABSTRACT

Graph data modeling and querying arises in many practical application domains such as social and biological networks where the primary focus is on concepts and their relationships and the rich patterns in these complex webs of interconnectivity. In this book, we present a concise unified view on the basic challenges which arise over the complete life cycle of formulating and processing queries on graph databases. To that purpose, we present all major concepts relevant to this life cycle, formulated in terms of a common and unifying ground: the property graph data model—the pre-dominant data model adopted by modern graph database systems.

We aim especially to give a coherent and in-depth perspective on current graph querying and an outlook for future developments. Our presentation is self-contained, covering the relevant topics from: graph data models, graph query languages and graph query specification, graph constraints, and graph query processing. We conclude by indicating major open research challenges towards the next generation of graph data management systems.

## KEYWORDS

# Contents

# Foreword

The current surge of interest in Graph Data Bases (GDBs) reflects the popularity of their data models based on nodes and edges, which, in many applications, provide a more intuitive conceptualization for entities and relationships than the one offered by Relational Data Bases (RDBs). This has inspired the design and development of many GDB systems and their use in a wide range of applications. Indeed to date, we counted more than 20 GDB systems developed and used in application areas such as Semantic Web, Social Networking, Fraud Detection, Recommendation Systems, Life Science, and Knowledge Bases.

For all their remarkable achievements, GDBs still lack the conceptual coherence that RDBs have been blessed with from the beginning as a result of E.F. Codd's seminal contributions which, combined with the major research advances in theory and systems that followed, provide the subject of numerous textbooks. However, the fast-expanding technology of GDBs is still quite far from achieving similar levels of conceptual unification and this create hurdles for researchers, instructors, and students alike.

This book tackles this problem head on by presenting a comprehensive unified treatment of GDBs, as needed to serve as a reference book for experts and a textbook for graduate students. The book's coverage begins with a formal treatment of the Property Graph Data Model that is common to most GDBs. Then, the book discusses GDB query languages and, moving past their many differences, it proposes a core property graph query language and elucidates its properties both in terms of graph logic and graph algebra. After that, the book covers techniques for efficient GDB implementation, including data structures, indexes, query operators, and processing, for which the presentation underscores how solutions different from those of traditional DBs are often required. Furthermore, the departures from traditional technology are even more dramatic for (i) integrity constraints, which lose their key role in normal-form RDB schema design, but find new important uses in GDBs, and (ii) interactive query specification via examples and counter-examples that have proven to be surprisingly effective with GDBs. The book's comprehensive treatment is further enhanced by extensive references and suggestions on open research problems for further investigation.

Carlo Zaniolo
Computer Science Department
University of California at Los Angeles (UCLA)

# Acknowledgments

The authors would like to warmly thank the many people who helped us to make this book a reality. First and foremost, we thank our families and partners for their patience and support throughout the many months dedicated to the writing of this book.

We also give many thanks to the Editor H. V. Jagadish and the Founding Editor M. Tamer Özsu for the opportunity and encouragement to publish this book. During the writing, the staff at Morgan & Claypool were just awesome, especially Diane Cerra. Thank you all for keeping the writing moving forward. We also thank the three reviewers for their critical and insightful feedback.

Our sincere thanks further go to Carlo Zaniolo for kindly writing the Foreword. We are greatly honored by your contribution!

Finally, we give our heartfelt thanks to colleagues for reading early drafts. We especially thank Sourav Bhowmick, Stefania Dumbrava, Jan Hidders, Wilco van Leeuwen, Davide Mottin, Oskar van Rest, and Kaijie Zhu for their detailed proofreading and helpful comments.

Angela Bonifati, George Fletcher, Hannes Voigt, and Nikolay Yakovets
September 2018