# Answering Queries Using Views

**Second Edition**

# Synthesis Lectures on Data Management

Answering Queries Using Views, Second Edition
Foto Afrati and Rada Chirkova
2019

Data Exploration Using Example-Based Methods
Matteo Lissandrini, Davide Mottin, Themis Palpanas, and Yannis Velegrakis
2018

Data Profiling
Ziawasch Abedjan, Lukasz Golab, Felix Naumann, and Thorsten Papenbrock
2018

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

Answering Queries Using Views, Second Edition

Foto Afrati and Rada Chirkova

# Answering Queries Using Views

## Second Edition

Foto Afrati
National Technical University of Athens

Rada Chirkova
North Carolina State University

*SYNTHESIS LECTURES ON DATA MANAGEMENT #54*

# ABSTRACT

The topic of using views to answer queries has been popular for a few decades now, as it cuts across domains such as query optimization, information integration, data warehousing, website design and, recently, database-as-a-service and data placement in cloud systems.

This book assembles foundational work on answering queries using views in a self-contained manner, with an effort to choose material that constitutes the backbone of the research. It presents efficient algorithms and covers the following problems: query containment; rewriting queries using views in various logical languages; equivalent rewritings and maximally contained rewritings; and computing certain answers in the data-integration and data-exchange settings. Query languages that are considered are fragments of SQL, in particular select-project-join queries, also called conjunctive queries (with or without arithmetic comparisons or negation), and aggregate SQL queries.

This second edition includes two new chapters that refer to tree-like data and respective query languages. Chapter 8 presents the data model for XML documents and the XPath query language, and Chapter 9 provides a theoretical presentation of tree-like data model and query language where the tuples of a relation share a tree-structured schema for that relation and the query language is a dialect of SQL with evaluation techniques appropriately modified to fit the richer schema.

# KEYWORDS

# Contents

# Preface to the First Edition

Views are used in various scenarios; some of the view-based settings have been considered in depth by the research community. The settings that we cover in this book include the following.

1. The problem of rewriting queries using views, where a set of views and a set of queries are given, and we need to find equivalent rewritings (if they exist) of the queries using these views.

2. Sometimes we cannot find equivalent rewritings but can still compute a significant part of the answer to the query. This gives rise to the problems of computing certain answers and of finding maximally contained query rewritings.

3. The picture of finding rewritings changes when we assume that the data satisfy certain constraints (dependencies). We re-examine the problem of finding rewritings for the setting in which the constraints are tuple-generating dependencies and equality-generating dependencies.

4. A closely related topic based on the same theoretical foundations is the data-exchange setting. We define the concept of certain answers and present algorithms to find them.

5. Some theoretical aspects of the more general problem of answering queries using views have also been investigated in more general and abstract settings, such as determining the (query) language fragments for which there are rewritings in the cases where the queries are determined by the views.

   In order to solve problems in the above settings, we need technical tools. Thus, we provide detailed treatments of some tools, including conjunctive-query containment (with and without arithmetic comparisons and negation), the chase algorithm for reasoning about dependencies, and going beyond nonrecursive languages to find rewritings, with a discussion of Datalog.

   What the book is not about: This book is not about indexes or data structures that implement the techniques considered in the exposition. Instead, this books focuses on the formal-logic perspective on the topic.

   The book is written in a linear way, in the sense that each chapter depends on all the previous chapters. We have made every effort for the book to be self-contained, thus there are no substantial prerequisites. A reader familiar with the basics of the theory of database systems and knowledge of logic will move faster through the chapters. Exercises are included. More exercises,

including online exercises, and any supplementary material will be found on the website[1] of the book.

Foto Afrati and Rada Chirkova
October 2017

---

[1] bit.ly/2zMYBHf

# Preface to the Second Edition

We have added two chapters that refer to tree-like data and respective query languages. Chapter 8 presents the data model for XML documents and the XPath query language and gives query containment tests and query rewriting techniques. However, besides the relational data model and the tree-structured XML data model, many recently developed systems, to manage big data, use data models that combine relational and tree-like features. In terms of defining and using views, these models have not be investigated. However, they are mature enough to be presented in a rigorous way. Thus, Chapter 9 provides a theoretical presentation of such data model and query language as extensions of the relational model and SQL query language. In these data models, the tuples of a relation share a tree-structured schema for that relation and the query language is a dialect of SQL. The query language uses SQL-style syntax but the evaluation techniques have to be modified to fit the richer schema. The conclusions and bibliographical notes for these two chapters are included in the chapters themselves. Chapters 8 and 9 can also be read independently of the rest of the chapters in the book and independently of each other.

In the chapters of the first edition, we have corrected errors.


Foto Afrati and Rada Chirkova
January 2019

# Acknowledgments