

# **Similarity Joins in Relational Database Systems**

# Synthesis Lectures on Data Management

Editor

**M. Tamer Özsu**, *University of Waterloo*

Synthesis Lectures on Data Management is edited by Tamer Özsu of the University of Waterloo. The series publishes 50- to 125 page publications on topics pertaining to data management. The scope will largely follow the purview of premier information and computer science conferences, such as ACM SIGMOD, VLDB, ICDE, PODS, ICDT, and ACM KDD. Potential topics include, but not are limited to: query languages, database system architectures, transaction management, data warehousing, XML and databases, data stream systems, wide scale data distribution, multimedia data management, data mining, and related subjects.

Similarity Joins in Relational Database Systems

Nikolaus Augsten and Michael H. Böhlen

2013

Data Cleaning: A Practical Perspective

Venkatesh Ganti and Anish Das Sarma

2013

Data Processing on FPGAs

Jens Teubner and Louis Woods

2013

Perspectives on Business Intelligence

Raymond T. Ng, Patricia C. Arocena, Denilson Barbosa, Giuseppe Carenini, Luiz Gomes, Jr.

Stephan Jou, Rock Anthony Leung, Evangelos Milios, Renée J. Miller, John Mylopoulos, Rachel A.

Pottinger, Frank Tompa, and Eric Yu

2013

Semantics Empowered Web 3.0: Managing Enterprise, Social, Sensor, and Cloud-based Data and Services for Advanced Applications

Amit Sheth and Krishnaprasad Thirunarayan

2012

Data Management in the Cloud: Challenges and Opportunities

Divyakant Agrawal, Sudipto Das, and Amr El Abbadi

2012

Query Processing over Uncertain Databases  
Lei Chen and Xiang Lian  
2012

Foundations of Data Quality Management  
Wenfei Fan and Floris Geerts  
2012

Incomplete Data and Data Dependencies in Relational Databases  
Sergio Greco, Cristian Molinaro, and Francesca Spezzano  
2012

Business Processes: A Database Perspective  
Daniel Deutch and Tova Milo  
2012

Data Protection from Insider Threats  
Elisa Bertino  
2012

Deep Web Query Interface Understanding and Integration  
Eduard C. Dragut, Weiyi Meng, and Clement T. Yu  
2012

P2P Techniques for Decentralized Applications  
Esther Pacitti, Reza Akbarinia, and Manal El-Dick  
2012

Query Answer Authentication  
HweeHwa Pang and Kian-Lee Tan  
2012

Declarative Networking  
Boon Thau Loo and Wenchao Zhou  
2012

Full-Text (Substring) Indexes in External Memory  
Marina Barsky, Ulrike Stege, and Alex Thomo  
2011

Spatial Data Management  
Nikos Mamoulis  
2011

Database Repairing and Consistent Query Answering  
Leopoldo Bertossi  
2011

Managing Event Information: Modeling, Retrieval, and Applications

Amarnath Gupta and Ramesh Jain

2011

Fundamentals of Physical Design and Query Compilation

David Toman and Grant Weddell

2011

Methods for Mining and Summarizing Text Conversations

Giuseppe Carenini, Gabriel Murray, and Raymond Ng

2011

Probabilistic Databases

Dan Suciu, Dan Olteanu, Christopher Ré, and Christoph Koch

2011

Peer-to-Peer Data Management

Karl Aberer

2011

Probabilistic Ranking Techniques in Relational Databases

Ihab F. Ilyas and Mohamed A. Soliman

2011

Uncertain Schema Matching

Avigdor Gal

2011

Fundamentals of Object Databases: Object-Oriented and Object-Relational Design

Suzanne W. Dietrich and Susan D. Urban

2010

Advanced Metasearch Engine Technology

Weiyi Meng and Clement T. Yu

2010

Web Page Recommendation Models: Theory and Algorithms

Sule Gündüz-Ögüdücü

2010

Multidimensional Databases and Data Warehousing

Christian S. Jensen, Torben Bach Pedersen, and Christian Thomsen

2010

Database Replication

Bettina Kemme, Ricardo Jimenez-Peris, and Marta Patino-Martinez

2010

### Relational and XML Data Exchange

Marcelo Arenas, Pablo Barcelo, Leonid Libkin, and Filip Murlak  
2010

### User-Centered Data Management

Tiziana Catarci, Alan Dix, Stephen Kimani, and Giuseppe Santucci  
2010

### Data Stream Management

Lukasz Golab and M. Tamer Özsu  
2010

### Access Control in Data Management Systems

Elena Ferrari  
2010

### An Introduction to Duplicate Detection

Felix Naumann and Melanie Herschel  
2010

### Privacy-Preserving Data Publishing: An Overview

Raymond Chi-Wing Wong and Ada Wai-Chee Fu  
2010

### Keyword Search in Databases

Jeffrey Xu Yu, Lu Qin, and Lijun Chang  
2009

© Springer Nature Switzerland AG 2022

Reprint of original edition © Morgan & Claypool 2014

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopy, recording, or any other except for brief quotations in printed reviews, without the prior permission of the publisher.

Similarity Joins in Relational Database Systems

Nikolaus Augsten and Michael H. Böhlen

ISBN: 978-3-031-00723-1      paperback

ISBN: 978-3-031-01851-0      ebook

DOI 10.1007/978-3-031-01851-0

A Publication in the Springer series

*SYNTHESIS LECTURES ON DATA MANAGEMENT*

Lecture #38

Series Editor: M. Tamer Özsu, *University of Waterloo*

Series ISSN

Synthesis Lectures on Data Management

Print 2153-5418    Electronic 2153-5426

# Similarity Joins in Relational Database Systems

Nikolaus Augsten  
University of Salzburg

Michael H. Böhlen  
University of Zürich

*SYNTHESIS LECTURES ON DATA MANAGEMENT #38*

## ABSTRACT

State-of-the-art database systems manage and process a variety of complex objects, including strings and trees. For such objects equality comparisons are often not meaningful and must be replaced by similarity comparisons. This book describes the concepts and techniques to incorporate similarity into database systems. We start out by discussing the properties of strings and trees, and identify the edit distance as the de facto standard for comparing complex objects. Since the edit distance is computationally expensive, token-based distances have been introduced to speed up edit distance computations. The basic idea is to decompose complex objects into sets of tokens that can be compared efficiently. Token-based distances are used to compute an approximation of the edit distance and prune expensive edit distance calculations.

A key observation when computing similarity joins is that many of the object pairs, for which the similarity is computed, are very different from each other. Filters exploit this property to improve the performance of similarity joins. A filter preprocesses the input data sets and produces a set of candidate pairs. The distance function is evaluated on the candidate pairs only. We describe the essential query processing techniques for filters based on lower and upper bounds. For token equality joins we describe prefix, size, positional and partitioning filters, which can be used to avoid the computation of small intersections that are not needed since the similarity would be too low.

## KEYWORDS

strings, trees, similarity, edit distance, q-grams, pq-grams, token-based distance, lower bound, upper bound, similarity join



*To Leni, Magdalena, and Katharina.*

*Nikolaus*

*To Franziska, Chantal, and Pascal.*

*Michael*

# Contents

<b>Preface</b> .....	<b>xv</b>
<b>Acknowledgments</b> .....	<b>xvii</b>
<b>1 Introduction</b> .....	<b>1</b>
1.1 Applications of Similarity Queries .....	1
1.2 Edit-Based Similarity Measures .....	4
1.3 Token-Based Similarity Measures .....	5
<b>2 Data Types</b> .....	<b>7</b>
2.1 Strings .....	7
2.2 Trees .....	7
<b>3 Edit-Based Distances</b> .....	<b>11</b>
3.1 String Edit Distance .....	11
3.1.1 Definition of the String Edit Distance .....	11
3.1.2 Computation of the String Edit Distance .....	13
3.2 Tree Edit Distance .....	15
3.2.1 Definition of the Tree Edit Distance .....	15
3.2.2 Computation of the Tree Edit Distance .....	18
3.2.3 Constrained Tree Edit Distance .....	18
3.2.4 Unordered Tree Edit Distance .....	19
3.3 Further Readings .....	22
<b>4 Token-Based Distances</b> .....	<b>25</b>
4.1 Sets and Bags .....	25
4.1.1 Counting Approach .....	25
4.1.2 Frequency Approach .....	25
4.2 Similarity Measures for Sets and Bags .....	26
4.2.1 Overlap Similarity .....	26
4.2.2 Jaccard Similarity .....	26
4.2.3 Dice Similarity .....	27

4.2.4	Converting Threshold Constraints	27
4.3	String Tokens	28
4.3.1	q-Gram Tokens	28
4.4	Tokens for Ordered Trees	29
4.4.1	Overview of Ordered Tree Tokens	30
4.4.2	The pq-Gram Distance	32
4.4.3	An Algorithm for the pq-Gram Index	37
4.4.4	Relational Implementation	38
4.5	Tokens for Unordered Trees	41
4.5.1	Overview of Unordered Tree Tokens	42
4.5.2	Desired Properties for Unordered Tree Decompositions	43
4.5.3	The Windowed pq-Gram Distance	46
4.5.4	Properties of Windowed pq-grams	50
4.5.5	Building the Windowed pq-Gram Index	56
4.6	Discussion: Properties of Tree Tokens	57
4.7	Further Readings	59
<b>5</b>	<b>Query Processing Techniques</b>	<b>61</b>
5.1	Filters	61
5.2	Lower and Upper Bounds	62
5.3	String Distance Bounds	63
5.3.1	Length Filter	63
5.3.2	Count Filter	63
5.3.3	Positional Count Filter	65
5.3.4	Using String Filters in a Relational Database	65
5.4	Tree Distance Bounds	69
5.4.1	Size Lower Bound	69
5.4.2	Intersection Lower Bound	69
5.4.3	Traversal String Lower Bound	70
5.4.4	pq-Gram Lower Bound	72
5.4.5	Binary Branch Lower Bound	74
5.4.6	Constrained Edit Distance Upper Bound	76
5.5	Further Readings	78
<b>6</b>	<b>Filters for Token Equality Joins</b>	<b>79</b>
6.1	Token Equality Join – Avoiding Empty Intersections	79
6.2	Prefix Filter – Avoiding Small Intersections	83

6.2.1	Prefix Filter for Overlap Similarity .....	83
6.2.2	Prefix Filter for Jaccard Similarity .....	85
6.2.3	Effectiveness of Prefix Filtering .....	86
6.3	Size Filter .....	87
6.4	Positional Filter .....	87
6.5	Partitioning Filter .....	88
6.6	Further Readings .....	88
<b>7</b>	<b>Conclusion .....</b>	<b>91</b>
	<b>Bibliography .....</b>	<b>93</b>
	<b>Authors' Biographies .....</b>	<b>103</b>
	<b>Index .....</b>	<b>105</b>

# Preface

During the last few decades database systems have evolved substantially and today it is common for database systems to manage a large variety of complex objects. At the physical level, there has been significant progress in terms of storing and processing such complex objects. At the logical level, however, complex objects remain a challenge. A key reason is that equality, which is appropriate for simple objects, is often ineffective for complex objects. In this book we describe the essential concepts and techniques toward a principled solution for processing complex objects that must be compared in terms of similarity rather than equality.

An intuitive approach to define the similarity of complex objects is the edit distance, i.e., the number of basic edit operations that are required to transform one object into another. The intuitive nature of the edit distance is the reason why edit distances have become the de facto standard for complex objects. We define the string and tree edit distance, give algorithms to compute these distances, and work out the essential properties that support the effective and efficient processing of complex objects.

Token distances, which decompose complex objects into sets of tokens, have been proposed to deal with the high computational cost of edit distances. The token distance is computed by comparing the token sets that are the result of the decompositions. The more similar two token sets are the smaller is their distance. Token sets are compared by counting the number of identical elements in the sets. Set intersection is an operation that is well supported by database systems and scales to large sets. We survey the different techniques to compute and process token sets, and we discuss in detail three representative decomposition techniques: strings with  $q$ -grams, ordered trees with  $pq$ -grams, and unordered trees with windowed  $pq$ -grams.

Determining the exact distance between complex objects, particularly for joins where all pairs of objects must be compared, is often too expensive. To reduce the costs of such computations, filter and refine approaches have been developed. The goal of the filter step is to cheaply identify candidate pairs for which the exact similarity must be computed. Non-candidate pairs do not have to be considered because their similarity is not sufficient to be included in the result. We describe various filter techniques, and provide lower and upper bounds that can be used to efficiently compute similarity joins.

The book uses strings and trees as representative examples of complex objects. The techniques discussed in this book, however, are general and are also applicable to other types of objects. In particular, graphs are an important data structure that recently have received a lot of attention, and for which edit- and token-based distances have been proposed. At the relevant places we provide references to the vibrant and emerging field of graphs in databases.

The book is intended as a starting point for researchers and students in the database field who would like to learn more about similarity in database systems. Throughout, we offer precisely defined concepts and properties, and we illustrate these with representative and carefully chosen examples. We made an effort to include precise definitions, theorems, and examples, but at the same time kept the description at a level that is understandable to a general audience with an academic background. Much of the material presented in this book has been used in courses taught during the last few years at the Free University of Bozen-Bolzano, Technische Universität München, University of Salzburg, and University of Zürich. Our warm thanks goes to the students who provided constructive feedback and helped to advance the material presented in this book.

Nikolaus Augsten and Michael H. Böhlen  
October 2013

# Acknowledgments

Several people offered valuable support during the preparation of this book. We warmly thank Tamer Özsu for inviting us to write this lecture, and we thank Diane Cerra for her perfect management of the entire process. We thank both Tamer and Diane for their kind but firm pushes that helped us to make this lecture happen.

Part of the book's material evolved from courses and lectures taught at the Free University of Bozen-Bolzano, Technische Universität München, the University of Salzburg, and the University of Zürich. We thank the many students for their constructive comments and their patience with early versions of the course material.

Tatsuya Akutsu and Robert Elsässer contributed with feedback about the unordered tree edit distance and its complexity, Willi Mann helped with proofreading the book. We would also like to acknowledge the many collaborators and friends who, through discussions and comments have shaped our thinking and understanding of the area: Arturas Mazeika, Chen Li, Johann Gamper, Denilson Barbosa, Themis Palpanas, Curtis Dyreson, Mateusz Pawlik, Jan Finis, Martin Raiber, Theo Härder, Leonardo Andrade Ribeiro, Benjamin Gufler, Gerard Lemson, Walter Costanzi, Franco Barducci, and Roberto Loperfido.

Finally, we would like to acknowledge our funding sources: Michael Böhlen's work is supported by the University of Zürich, the Swiss National Science Foundation, and the Free University of Bozen-Bolzano. Nikolaus Augsten's work is supported by the University of Salzburg, the Free University of Bozen-Bolzano, the Department for Promotion of Educational Policies, University and Research of the Autonomous Province of Bolzano - South Tyrol, and Technische Universität München.

Nikolaus Augsten and Michael H. Böhlen  
October 2013