

Chapter 11

Estimating the Completeness of Range Queries over Structured P2P Databases: Fundamentals, Theory, and Effective Applications to Distributed Information Systems

Alfredo Cuzzocrea

ICAR-CNR, Italy & University of Calabria, Italy

Marcel Karnstedt

DERI, NUI Galway, Ireland

Manfred Hauswirth

DERI, NUI Galway, Ireland

Kai-Uwe Sattler

Ilmenau University of Technology, Germany

Roman Schmidt

Ecole Polytechnique Federale de Lausanne, Switzerland

ABSTRACT

Range queries are a very powerful tool in a wide range of data management systems and are vital to a multitude of applications. The hierarchy of structured overlay systems can be utilized in order to provide efficient techniques for processing them, resulting in the support of applications and techniques based on range queries in large-scale distributed information systems. On the other hand, due to the rapid development of the Web, applications based on the P2P paradigm gain more and more interest, having such systems started to evolve towards adopting standard database functionalities in terms of complex query processing support. This goes far beyond simple key lookups, as provided by standard distributed hash tables (DHTs) systems, which makes estimating the completeness of query answers a crucial chal-

DOI: 10.4018/978-1-4666-1794-0.ch011

lenge. Unfortunately, due to the limited knowledge and the usually best-effort characteristics, deciding about the completeness of query results, e.g., getting an idea when a query is finished or what amount of results is still missing, is very challenging. There is not only an urgent need to provide this information to the user issuing queries, but also for implementing sophisticated and efficient processing techniques based on them. In this chapter, the authors propose a method for solving this task. They discuss the applicability and quality of the estimations, present an implementation and evaluation for the P-Grid system, and show how to adapt the technique to other overlays. The authors also discuss the semantics of completeness for complex queries in P2P database systems and propose methods based on the notion of routing graphs for estimating the number of expected query answers. Finally, they discuss probabilistic guarantees for the estimated values and evaluate the proposed methods through an implemented system.

INTRODUCTION

Many new applications on the Web are based on the idea of collecting and combining large public data sets and services. In such public data management applications, the information, its structure and its semantics in many cases are the result of the collaborative effort of the participants. Examples of such applications are social networks, e.g., friend-of-a-friend networks, distributed recommender systems, distributed directory and index services, and sharing of sensor data. These applications typically require the indexing and management of data distributed over a large number of independent data stores, which is a typical scenario targeted by overlay networks.

P2P systems and particularly structured overlays based on distributed hashtables (DHTs) are recognized as a promising infrastructure for large-scale distributed data management. The main reasons are their effectiveness and scalability as well as the predictable behavior. After the first generation supporting only basic key lookups, recent research efforts address also the problem of efficiently querying range predicates (Bharambe, A., Agrawal, M., Seshan, S. (2004); Datta, A., Hauswirth, M., Schmidt, R., John, R., Aberer, K. (2005)). Typically, these approaches exploit the structure of the overlay (e.g., a tree structure) by implementing some kind of multicast protocol. In this context, a main challenge is to

estimate the progress of query processing, i.e., to answer the question which fraction of the total query result is already received. The difficulties are due to the purely decentralized nature of the structured overlay, the lack of global knowledge (no peer knows how many peers are responsible for the queried key range), the dynamics of the network (peers may leave the network during processing a query), as well as the often used best-effort strategy for query routing and answering. Indeed, P2P data management is inherently open world: While processing a query, peers can fail, leave or join the network, or simply send no or a delayed answer (Gribble, S.D., Halevy, A.Y., Ives, Z.G., Rodrig, M., Suciu, D. (2001)). Though this can be mitigated by replication and delay-tolerant query techniques, there is no guarantee that all answers which potentially exist can be returned. On the other hand, DHTs, the most efficient family of overlay networks, so far have only been applicable to a certain degree in these scenarios, as support for managing and querying structured data in DHTs still is limited.

Estimating the completeness of a query result is not only a helpful information for the user issuing the query, but it is also needed for processing complex queries. For instance, query operators like aggregation or ranking-based queries (e.g., skyline queries (Börzsönyi, S., Kossman, D., Stocker, K. (2001); Karnstedt, M., Müller, J., Sattler, K. (2007))) require to know when all input data is

29 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/estimating-completeness-range-queries-over/67001

Related Content

Autonomic Networking Integrated Model and Approach (ANIMA): Secure Autonomic Network Infrastructure

Toerless Eckert (2019). *Emerging Automation Techniques for the Future Internet* (pp. 90-112).

www.irma-international.org/chapter/autonomic-networking-integrated-model-and-approach-anima/214428

Analyzing the Capacity of Unsolicited Political Email

Kristin Johnson and Brian S. Krueger (2012). *E-Politics and Organizational Implications of the Internet: Power, Influence, and Social Change* (pp. 220-244).

www.irma-international.org/chapter/analyzing-capacity-unsolicited-political-email/65217

Aspect-Oriented Programming and Aspect.NET as Security and Privacy Tool for Web and 3D Web Programming

Vladimir O. Safonov (2011). *Security in Virtual Worlds, 3D Webs, and Immersive Environments: Models for Development, Interaction, and Management* (pp. 221-262).

www.irma-international.org/chapter/aspect-oriented-programming-aspect-net/49524

Integrating Big Data to Smart Destination Heritage Management

Kubra Ozer, Mehmet Altug Sahin and Gurel Cetin (2022). *Handbook of Research on Digital Communications, Internet of Things, and the Future of Cultural Tourism* (pp. 411-429).

www.irma-international.org/chapter/integrating-big-data-to-smart-destination-heritage-management/295515

Data Caching in Web Applications

Tony C. Shan and Winnie W. Hua (2008). *Encyclopedia of Internet Technologies and Applications* (pp. 132-141).

www.irma-international.org/chapter/data-caching-web-applications/16845