Lecture Notes in Computer Science 3128

Commenced Publication in 1973 Founding and Former Series Editors: Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Editorial Board

Takeo Kanade
Carnegie Mellon University, Pittsburgh, PA, USA
Josef Kittler
University of Surrey, Guildford, UK
Jon M. Kleinberg
Cornell University, Ithaca, NY, USA
Friedemann Mattern
ETH Zurich, Switzerland
John C. Mitchell
Stanford University, CA, USA
Moni Naor
Weizmann Institute of Science, Rehovot, Israel
Uscar Nierstrasz
C Den du Den con
Indian Institute of Technology Madras India
Bernhard Steffen
University of Dortmund, Germany
Madhu Sudan
Massachusetts Institute of Technology, MA, USA
Demetri Terzopoulos
New York University, NY, USA
Doug Tygar
University of California, Berkeley, CA, USA
Moshe Y. Vardi
Rice University, Houston, TX, USA
Gerhard Weikum
Max-Planck Institute of Computer Science, Saarbruecken, Germany

Dmitri Asonov

Querying Databases Privately

A New Approach to Private Information Retrieval



Author

Dmitri Asonov IBM Almaden Research Center 650 Harry Road, San Jose, CA 95123, USA E-mail: dasonov@us.ibm.com

Dissertation der Humboldt Universität zu Berlin Tag der mündlichen Prüfung: 7. Juli 2003

Referent: Prof. Johann-Christoph Freytag, Ph.D., Humboldt Universität zu Berlin Referent: Prof. Oliver Günther, Ph.D., Humboldt Universität zu Berlin Referent: Rakesh Agrawal, Ph.D., IBM Almaden Research Center

Library of Congress Control Number: 2004094685

CR Subject Classification (1998): H.3, H.2, H.4, K.4, K.6.5, C.2

ISSN 0302-9743 ISBN 3-540-22441-6 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable to prosecution under the German Copyright Law.

Springer-Verlag is a part of Springer Science+Business Media

springeronline.com

© Springer-Verlag Berlin Heidelberg 2004 Printed in Germany

Typesetting: Camera-ready by author, data conversion by Olgun Computergrafik Printed on acid-free paper SPIN: 11018001 06/3142 5 4 3 2 1 0

Foreword

The Internet and the World Wide Web (WWW) play an increasingly important role in our today's activities. More and more we use the Web to buy goods and to inform ourselves about cultural, political, economical, medical, and scientific developments. For example, accessing flight schedules, medical data, or retrieving stock information become common practice in today's world. Many people assume that there is no one who "watches" them when accessing this data.

However, sensitive users who access electronic shops (e-shops) might have observed that this assumption often is not true. In many cases, E-shops track the users' "access behavior" when browsing the Web pages of the e-shop thus deriving "access patterns" for individual shoppers. Therefore, this knowledge on access behavior and access patters allows the system to tailor access to Web pages for that user to his/her specific needs in the future. This tracking of users might be considered harmless and "acceptable" in many cases. However, in cases when this information is used to harm a person - for example about the person's health problems - or to violate his/her privacy (for example finding out about his/her financial situation), he/she would like to be sure that such tracking is impossible to protect the user's rights.

These simple examples clearly demonstrate the necessity to shield the user from such spying to protect his/her privacy. That is, a user should be able to access a database (or a data source in general) without allowing others to "observe" which data is requested and accessed by the user; neither the query nor the answer should be visible or accessible to others. Surprisingly, despite the urgent need for concepts and techniques to protect the user from being spied on, very few results are known and available that addresses the problem adequately. During the last 10 years the area of **Private Information Retrieval (PIR)** has addressed some of the problems concerning privacy. However many of those results are of theoretical nature and thus do not carry over into practical solutions for protecting privacy when accessing information sources on the Web or in databases.

With this book Dr. Asonov is one of the first researchers who addresses the topic of querying data privately in a systematic and comprehensive way developing practical solutions in the context of database systems. The results presented in this book sometimes might look theoretical, but they describe his clear understanding of the problem as well as the solutions required for "real-world" settings, in particular for scalable database solutions. As a basis Dr. Asonov first presents the framework for privately accessing database by developing several algorithms which also include the use of special hardware. In the second part of the book he focuses on solving several important subproblems; for them he also includes some validation by benchmarking to show to efficiency of the solutions. Finally, Dr. Asonov shows how his solutions could be used in solving some problems in the area of voting and digital rights management. Initially, these problems seem to be completely unrelated to PIR, however Dr. Asonov shows how some of his results can be used for creative solutions in the areas mentioned. Overall, the careful reader will notice that - despite the many technical details -his in-depth treatment of privacy in database provides the insight into the problem necessary for such an important topic.

In summary, with this book Dr. Asonov provides a systematic treatment of the problem how to access databases privately. The way he approaches the problem and how he develops solutions makes this book valuable for both researchers and practitioners who are interested in better understanding the issues. He develops scalable solutions that are necessary and important in the context of private information retrieval/private database access. The in-depth presentation of the algorithms and techniques is enlightening to students and a valuable resource for computer scientists. I predict that this book will provide the "starting point" for others to perform further research and development in this area.

Prof. Johann-Christoph Freytag, Ph.D., May 2004

Preface

People often retrieve information by querying databases. Designing databases that allow a user to execute queries *efficiently* is a subject that has been investigated for decades, and is now often regarded as a "researched to death" topic. However, the evolution of information technologies and society makes the database area a consistent source of new, previously unimaginable research challenges. This work is dedicated to partially meeting one of these new challenges: querying databases *privately*.

This new challenge is due to a very fundamental constraint of the conventional concept of querying information. Namely, in the conventional setting, the one who queries (the user) must reveal the query content and, by implication, the result of querying to the one who processes the query (the database server). This constraint seems to be negligible if the user trusts the server. However, the growing population of information providers makes it extremely difficult for users to establish and rely on the trustworthiness of information providers. Indeed, more and more cases are reported wherein information providers misuse the information provided by users' queries against the users, for example by sharing this information with third parties without permission, or by using this information for unsolicited advertisement.

We approach this constraint in a direct manner: If it is difficult to trust the server, we could try to remove the need for trust completely, by hiding the content of the user query and result from the server. This research problem, called Private Information Retrieval (PIR), has been under intensive and mainly theoretical investigation since 1996. These results are classified and analyzed in the first of four parts of this book. Our main contribution is considering this problem from a practical angle, as follows.

In Part II, we accept the assumptions and simplifications made in previous related work, and focus on obtaining efficient solutions and algorithms without changing the common model. Namely, we break the established belief that the server must read the entire database for a PIR protocol to answer a query. We further develop our solution by improving the processing and preprocessing complexities of our PIR protocol.

In Part III we extend the common PIR model in two directions. First, we relax the requirement that no information about a query must be revealed. This allows us to offer the user a trade-off between the level of privacy required and the response time for a query. The second extension of the model is done by understanding the economics associated with the PIR problem. Namely, we assumed that information in the database is from different owners. We then consider the problem of distributing royalties between the information owners, given that no information about content of user queries is revealed.

A number of questions remain to be answered before the problem of querying databases privately can be regarded as completely investigated. However, we argue that results presented in the book have pushed the state of the art in this area, from the entirely theoretical level to the stage where implementing an applicable prototype can be considered ultimately possible.

Acknowledgements

I am most indebted to Professor Johann-Christoph Freytag for the success of this work. Our interaction was an example of a brilliant collaboration between a student and an adviser, so rarely found in science.

I was lucky to secure Professor Oliver Günther as my second advisor. I learned a lot from him. Professor Günther naturally supplemented the image of a perfect professor that I perceived from my first advisor.

I am very grateful to Rakesh Agrawal from IBM Almaden Research Center for being an external reviewer of my dissertation. Professor Sean W. Smith and Alex Iliev from Dartmouth College, Ronald Perez from IBM T. J. Watson Research Center, Christian Cachin from IBM Zürich Research Laboratory, and Frank Leymann from IBM Laboratory Böblingen were my occasional, but nevertheless most valuable external contacts.

I could not survive the hardship of making a Ph.D. without the warm, social support from my graduate school colleagues, and the team of DBIS department of Humboldt University. Especially, I would like to thank Markus Schaal and Christoph Hartwich for our fruitful collaboration in CS research, and my officemates Felix Naumann and Heiko Müller, who had to listen to my erroneous German every day. Ulrike Scholz and Heinz Werner have made DBIS a very comfortable place to work at.

My russian–speaking friends in Berlin, Stanislav Isaenko, Viktor Malyarchuk, and Mykhaylo Semtsiv helped me better understand research as a process by sharing their experiences in biological and physical research.

My teachers in Moscow provided the educational background from which I am benefiting now. Among them Yulia A. Azovzeva, Alexei I. Belousov, Valeri M. Chernenki, Maria T. Lepeshkina, Sergei V. Nesterov, Valentina P. Strekalova, Sergei A. Trofimov, and Valeri D. Vurdov were most helpful.

Last but not least, I am thankful to my family who supported me all the way through.

This research was supported by the German Research Society, Berlin-Brandenburg Graduate School in Distributed Information Systems (DFG grants no. GRK 316 and GRK 316/2).

Table of Contents

Part I. Introduction and Related Work

1	T 4 .		9
T	Inti		3
	1.1	Problem Statement	- 3
	1.2	Book Outline	6
	1.3	Motivating Examples	8
		1.3.1 Examples of Violation of User Privacy	8
		1.3.2 Application Areas for PIR	9
2	\mathbf{Rel}	ated Work	11
	2.1	Naive Approaches Do Not Work	11
	2.2	PIR Approaches	11
		2.2.1 Theoretical Private Information Retrieval	12
		2.2.2 Computational Private Information Retrieval	13
		2.2.3 Symmetrical Private Information Retrieval	14
		2.2.4 Hardware-Based Private Information Retrieval	14
		2.2.5 Further Extensions of the Problem Setting	16
		2.2.6 PIR with Preprocessing and Offline Communication	17
		2.2.7 Work Related to PIR Indirectly	18
	2.3	Analysis of the Previous Approaches	18
		2.3.1 Evaluation Criteria for PIR Approaches	18
		2.3.2 State of the Art	19
		2.3.3 Open Problems	20

Part II. Almost Optimal PIR

3	PIR wit	th $O(1)$ Query Response Time	
	and $O(1)$	1) Communication	3
	3.1 Bas	ic Protocol 2	23
	3.1.	1 Database Shuffling Algorithm (SSA) 2	24
	3.1.	2 The Protocol	26
	3.1.	3 An Algorithm for Processing a Query 2	27
	3.1.	4 Trade-Off between Preprocessing Workload	
		and Query Response Time 2	27

		3.1.5 Choosing the Optimal Trade-Off	28
		3.1.6 Multiple Queries and Multiple Coprocessors	30
	3.2	Formal Definition of the Privacy Property	30
		3.2.1 Basics of Information Theory	31
		3.2.2 Privacy Definition	33
	3.3	Proof of the Privacy Property of the Protocol	34
	3.4	Summary	35
4	Imp	proving Processing and Preprocessing Complexity	37
	4.1	Decreasing Query Response Time	37
	4.2	Decreasing the Complexity of Shuffling	38
		4.2.1 Split-Shuffle-Gather Algorithm (SSG)	38
		4.2.2 Balancing the Preprocessing Complexity between SC	
		and UC	41
		4.2.3 Recycling Used Shuffled Databases	42
	4.3	Measuring Complexity of the PIR Protocols	44
		4.3.1 A Normalized Measure for the Protocol Complexity	44
		4.3.2 The Measurement	45
	4.4	Summary	46
5	Exp	perimental Analysis of Shuffling Algorithms	49
	5.1	Shuffling Based on Bitonic Sort (SBS)	49
	5.2	Experiments	49
		5.2.1 Setup Details	50
		5.2.2 Experimental Data Collected	51
		5.2.3 Analysis	53
	5.3	The Superiority of SSG	53
		5.3.1 Imperfection of the Theoretically	
		Estimated Complexity of SSG	53
		5.3.2 On Minimal Bound for Shuffling Complexity	54
	5.4	Summary	55

Part III. Generalizing the PIR Model

6	Rep	oudiative Information Retrieval	59
6.1 The Need for Trade-Off between Privacy and Complexity			
		6.1.1 Our Results	60
		6.1.2 Preliminaries and Assumptions	60
	6.2 Defining Repudiation and Assessing Its Robustness		
		6.2.1 Repudiation Property	60
		6.2.2 Assessing the Robustness of Repudiation	62
6.3 Basic Repudiative Information Retrieval Protocol		Basic Repudiative Information Retrieval Protocol	64
		6.3.1 Analyzing the Robustness of the Protocol	65
		6.3.2 Multiple Queries	66

	6.3.3	Complexity of Preprocessing	68
	6.3.4	Summary of the Basic RIR	68
6.4	Varyi	ng the Robustness of the RIR Protocol	68
	6.4.1	A Parameterized RIR Protocol	69
	6.4.2	How Parameters Determine Robustness of Repudiation	69
	6.4.3	Turning the RIR Protocol into a PIR Protocol	71
6.5	Relate	ed Work	71
	6.5.1	Deniable Encryption	72
	6.5.2	Alternatives to the Quantification of Repudiation	72
6.6	Discus	ssion	73
	6.6.1	Redefining Repudiation	73
	6.6.2	Yet Another Alternative to the Quantification	
		of Repudiation	74
	6.6.3	Misinforming the Observers	74
6.7	Summ	nary	75
Dig	ital Ri	ights Management for PIR	77
Dig 7.1	ital Ri The C	ights Management for PIR	77 77
Dig 7.1 7.2	ital R i The C DRM	ights Management for PIR Collision between DRM and PIR without Repudiation	77 77 78
Dig 7.1 7.2 7.3	ital Ri The C DRM RIR S	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM	77 77 78 80
Dig 7.1 7.2 7.3 7.4	ital R The C DRM RIR S Robus	ights Management for PIRCollision between DRM and PIRwithout RepudiationSupporting DRMstness of Repudiation vs. Precision	77 77 78 80
Dig 7.1 7.2 7.3 7.4	ital Ri The C DRM RIR S Robus of Roy	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution	77 77 78 80 80
Dig 7.1 7.2 7.3 7.4 7.5	ital Ri The C DRM RIR S Robus of Roy The I	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Drawback of the Proposed DRM Scheme	77 77 78 80 80 81
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital Ri The C DRM RIR S Robus of Roy The I Absol	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Orawback of the Proposed DRM Scheme ute Privacy in Voting	77 77 78 80 80 81 84
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R The C DRM RIR S Robus of Roy The I Absol 7.6.1	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Orawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries	77 77 78 80 80 81 84 85
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Drawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions	77 77 78 80 80 81 84 85 88
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2 7.6.3	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Orawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions Probabilistic Voting Functions	77 77 78 80 80 81 84 85 88 90
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2 7.6.3 7.6.4	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Drawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions Probabilistic Voting Functions Related Work	77 77 78 80 81 84 85 88 90 93
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2 7.6.3 7.6.4 7.6.5	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision yalty Distribution Drawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions Probabilistic Voting Functions Related Work Discussion	77 77 78 80 81 84 85 88 90 93 95
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2 7.6.3 7.6.4 7.6.5 7.6.6	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision valty Distribution Drawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions Probabilistic Voting Functions Related Work Discussion The Implication of Absolute Privacy	77 77 78 80 81 84 85 88 90 93 95 96
Dig 7.1 7.2 7.3 7.4 7.5 7.6	ital R: The C DRM RIR S Robus of Roy The I Absol 7.6.1 7.6.2 7.6.3 7.6.4 7.6.5 7.6.6 Summ	ights Management for PIR Collision between DRM and PIR without Repudiation Supporting DRM stness of Repudiation vs. Precision valty Distribution Orawback of the Proposed DRM Scheme ute Privacy in Voting Preliminaries Deterministic Voting Functions Probabilistic Voting Functions Related Work Discussion The Implication of Absolute Privacy	77 77 78 80 81 84 85 88 90 93 95 96 96

Part IV. Discussion

 $\mathbf{7}$

8	Con	clusio	on and Future Work	101
	8.1	Summ	nary	101
	8.2	Future	e Work	104
		8.2.1	Querying Databases Privately without Tamper-Resist	ant
			Hardware	104
		8.2.2	Elaborate Query–Database Models	105
Ref	feren	ces		107
Ind	ex			115