

Micro-views, or on How to Protect Privacy while Enhancing Data Usability - Concepts and Challenges

Ji-Won Byun
byunj@cs.purdue.edu

Elisa Bertino
bertino@cerias.purdue.edu

CERIAS and Department of Computer Science
Purdue University
656 Oval Drive, West Lafayette, IN 47907

ABSTRACT

The large availability of repositories storing various types of information about individuals has raised serious privacy concerns over the past decade. Nonetheless, database technology is far from providing adequate solutions to this problem that requires a delicate balance between an individual's privacy and convenience and data usability by enterprises and organizations - a database which is rigid and over-protective may render data of little value. Though these goals may seem odd, we argue that the development of solutions able to reconcile them will be an important challenge to be addressed in the next few years. We believe that the next wave of database technology will be represented by a DBMS that provides high-assurance privacy and security. In this paper, we elaborate on such challenges. In particular, we argue that we need to provide different views of data at a very fine level of granularity; conventional view technology is able to select only up to a single attribute value for a single tuple. We need to go even beyond this level. That is, we need a mechanism by which even a single value inside a tuple's attribute may have different views; we refer them as micro-views. We believe that such a mechanism can be an important building block, together with other mechanisms and tools, of the next wave of database technology.

1. INTRODUCTION

Current information technology enables many organizations to collect, store and use a vast amount of personal information in their databases. The use of innovative knowledge extraction techniques combined with advanced data integration and correlation techniques [7, 8, 16] makes it possible to automatically extract a large body of information from the available databases and from a large variety of information repositories available on the web. Such a wealth of information and extracted knowledge raises, however, serious concerns about the privacy of individuals. As privacy awareness increases, individuals are becoming more reluctant to carry out their businesses and transactions online, and many enterprises are losing a considerable amount of potential profits [12]. Also, enterprises that collect information about individuals are in effect obligated to keep the collected information private and must strictly control the use of such information. Thus, information stored in the databases of an enterprise is not only a valuable property of the enterprise, but also a costly responsibility. Consequently, data management techniques providing high-assurance privacy and at the same time avoiding unnecessary restrictions to data access are in great demand. Equipped with such techniques, organizations shall be able to utilize information analysis and

knowledge extraction to provide better and tailored services to individuals without violating individual privacy.

To date, issues related to privacy have been widely investigated and several privacy protecting techniques have been developed. To our best knowledge, the most well known effort is the W3C's Platform for Privacy Preference (P3P) [20]. P3P allows websites to express their privacy policy in a machine readable format so that using a software agent, consumers can easily compare the published privacy policies against their privacy preferences. P3P, however, does not provide any functionality to keep these promises in the internal privacy practice of enterprises. To complement P3P's lack of enforcement mechanisms, many privacy-aware access control models have also been investigated [3, 4, 9, 10]. Although all these models do protect privacy of data providers¹, they are very rigid and do not provide ways to maximize the utilization of private information. Specifically, in those models access decision is always binary; i.e., a data access is either allowed or denied as in most conventional access control models.

We believe that a new generation of privacy-aware access control models should maximize information usability by exploiting the nature of information privacy. First, information privacy is context-specific. For instance, consider address data of consumers. The tolerance level of individuals for their address being used for direct marketing could be significantly different from the address being used for consumer analysis. Furthermore, the tolerance level varies from individual to individual. Some consumers may feel that it is acceptable to disclose their purchase history or browsing habits in return for better services; others may feel that revealing such information violates their privacy. These differences in individuals suggest that access control models should be able to maximize the utilization of private information by taking such large variations into account.

Second, the use of data generalization² can significantly increase the comfort level of data providers. For example, suppose that an enterprise has collected the birth dates of its consumers. Such information is very personal, and many individuals may not be comfortable with their information being used. Suppose now that the enterprise promises its consumers that this information will be used only in a generalized form; e.g., $\langle 07/23/1970 \rangle$ will be generalized to a less specific value $\langle 07/1970 \rangle$ or a categorical value $\langle 1965 - 1975 \rangle$. This assurance will surely comfort many consumers. Clearly,

¹By data providers, we refer to the subjects to whom the stored data is related.

²Data generalization refers to techniques that "replace a value with a less specific but semantically consistent value" [17].

Term	Description	Example
Privacy level	Level of privacy required by data provider	Low, Medium, High
Data type	Types of data being collected	Name, Address, Income, Age
Data usage type	Types of potential data usage (i.e. purpose)	Marketing, Admin, Shipping

Table 1: Privacy level, data type and data usage type

privacy enhancing access control models should be able to utilize more information by employing data generalization techniques.

The development of a DBMS that addresses the above requirements is a challenging task which requires revisiting theoretical foundations of data models as well as languages and architectures. A core DBMS component which is crucial in such a context is the access control system. Current access control systems are fundamentally inadequate with respect to the above goals. For example, fine-grained access control of data, an important requirement for privacy, poses several difficult problems and to date no satisfactory solution exists. We have yet to understand the relevant technical requirements.

In this paper, we pose as a new challenge the development of a new generation of access control systems. As an example, we propose a radically new access control model that is able to exploit the subtle nature of information privacy to maximize the usability of private information for enterprises with privacy guarantees. Our model is not to be considered a complete solution; rather it is meant to show some of capabilities that, in our opinion, a suitable model should provide. In particular, our model is based on the notion of a *micro-view*. A micro-view applies the well known idea of views at the level of the atomic components of tuples, that is, to an attribute value. By using the different precisions, one is able to finely calibrate the amount of information released by queries.

The remainder of this paper is organized as follows. In Section 2, we present a high-level description of our access control model. We discuss the technical challenges raised by our model in Section 3 and provide a brief survey of related work in Section 4. We then conclude our discussion in Section 5.

2. A SKETCH OF OUR “NAIVE” MODEL

Our model is based on a typical life-cycle of data concerning individuals. During the data collection phase, a data provider submits her privacy requirement, which specifies permissible usages of each data item and a level of privacy for each usage. This requirement is then stored in the database along with the collected data, and access to the data is strictly governed by the data provider’s requirement. In this section, we first illustrate how data collection is carried out in our model and discuss the access control model in detail.

2.1 Data collection and preprocess

As a prerequisite to our approach, supported privacy levels, types of data and possible data usages (i.e., purposes) have to be clearly defined and clarified through the published privacy policy. Table 1 describes these concepts with some practical examples.

For simplicity of discussion, we consider the following privacy levels only: *Low*, *Medium* and *High*. We also consider only three data types, *name*, *address* and *income*, and two

usage types, *admin* and *marketing*.

As previously mentioned, when individuals release their personal information, they specify permissible usages of each of their data items and a level of privacy for each usage. For instance, a data provider may select *Low* on *Address* for *Admin*; that is, she does not have any privacy concern over the address information when it is used for the purpose of administration. Thus, the address information can be used for the administrative purpose without any modification. However, the data provider may select *High* on *Address* for *marketing*. This indicates that she has great concerns about privacy of the address information when it is used for the purpose of marketing; thus, the address information should be used only in a sufficiently generalized form for the marketing purpose. Note that if one wishes to allow data providers to completely opt out from any use of data, another privacy level (e.g., Opt-out) can be added to indicate that the particular data should not be used at all.

Additional to storing the specified privacy requirements, the actual data items are preprocessed before being stored in the following way. Each data item is generalized and stored according to a multilevel organization, where each level corresponds to a specific privacy level. Intuitively, data for a higher privacy level requires a higher degree of generalization. For instance, the address data is stored into three levels: entire address for *Low*, city and state for *Medium* and state for *High*.

Table 2 illustrates some fictional records and privacy requirements stored in a conceptual database relation. Notice that every data item is stored in three different generalization levels, each of which corresponds to a particular privacy level. *PL_Admin* and *PL_Marketing* are metadata columns³ storing the set of privacy levels of data for *Admin* and *Marketing*, respectively. For instance, {L, L, M} in *PL_marketing* indicates that for the marketing purpose the privacy levels of *Name* and *Address* are both *Low* while the privacy level of *Income* is *Medium*.

Note that exactly how such data is organized and stored in the database is a crucial issue as it determines the performance and storage efficiency. However, this issue is beyond the scope of this paper and shall be discussed in our future work.

2.2 Access Control

In our model, data users query the database using standard SQL statements. However, the data accessible to each query varies depending on the privacy levels of the data and the purpose of the query⁴. That is, each query runs as if it is running on a view that is defined by the purpose of the query and the privacy levels of data. We call such views *micro-views*. Tables 3 and 4 illustrate this effect. For in-

³The metadata columns may be viewable to any user, but they can be modified only by authorized users.

⁴In this paper, we assume that each query is associated with a specific purpose.

CustID	Name		Address		Income		PL_Admin	PL_Marketing
1001	L	Alice Park	L	123 First St., Seattle, WA	L	45,000	{L, M, H}	{H, H, H}
	M	Alice P.	M	Seattle, WA	M	40K-60K		
	H	A.P.	H	WA	H	Under 100K		
1002	L	Aaron Parker	L	491 3rd St, Lafayette, IN	L	121,000	{L, L, M}	{H, M, H}
	M	Aaron P.	M	Lafayette, IN	M	120K-140K		
	H	A.P.	H	IN	H	Over 100K		
1003	L	Carol Jones	L	35 Oval Dr, Chicago, IL	L	64,000	{L, L, L}	{L, M, H}
	M	Carol J.	M	Chicago, IL	M	60K-80K		
	H	C.J.	H	IL	H	Under 100K		

Table 2: Private information and metadata

CustID	Name	Address	Income
1001	Alice Park	Seattle	Under 100K
1002	Aaron Parker	491 3rd St, Lafayette, IN	120K-140K
1003	Carol Jones	35 Oval Dr, Chicago, IL	64,000

Table 3: Micro-view for *Admin* purpose

CustID	Name	Address	Income
1001	A. P.	WA	Under 100K
1002	A. P.	Lafayette, IN	Over 100K
1003	Carol Jones	Chicago, IL	Under 100K

Table 4: Micro-view for *Marketing* purpose

stance, any query against the base table in Table 2 with *Admin* purpose returns a result that is equivalent to the result of the query run on the micro-view in Table 3. As the micro-views directly reflect the information that is allowed by each data provider, querying against these views does not violate privacy.

Note that the major difference of our model from conventional database models is that in our model, different sets of data may be returned for the same query, depending on the privacy levels of data and the purpose of the query. For instance, suppose that the following query is written against the base table in Table 2: “**SELECT * FROM Customer WHERE CustID = 1002**”. If the purpose of this query is *Admin*, then the system will return a tuple (‘Aaron Parker’, ‘491 3rd St, Lafayette, IN’, ‘120K-140K’) as Aaron’s privacy levels for *Admin* are specified as {L, L, M}. On the other hand, if the purpose of the query is *Marketing*, then a tuple (‘A. P.’, ‘Lafayette, IN’, ‘Over 100K’) will be retrieved as his privacy levels for *Marketing* is {H, M, H}.

An important issue to be addressed is how we associate a purpose with each query. Note that it is not trivial for a system to correctly infer the purpose of a query as it means that the system must correctly deduce the actual intention of database users. However, if we assume that users are honest, then the problem of associating a purpose with each query becomes relatively easy; i.e., users themselves can specify the purpose of their queries with an additional clause. For instance, a simple select statement “**SELECT name FROM customer**” can be extended to a form of “**SELECT name FROM customer FOR marketing**”. We believe that this is a reasonable approach. Many privacy violations occur from accidentally accessing unauthorized information, and thus it is important to develop a mechanism that database users can use to protect themselves from committing such accidental violations. A more sophisticated approach which validates whether users are indeed permitted to use their claimed purposes is thoroughly investigated in [5].

3. CHALLENGES

The comprehensive development of the approach we have sketched in the previous section and its integration in a DBMS architecture requires addressing several interesting challenges.

Policy specification language. The core of our model is that data providers can specify their privacy requirements using a privacy level for each data category. There is thus a strong need for a language in which privacy specifications can be expressed precisely. A challenge is that the language must be powerful enough to express every possible requirement, yet simple enough to avoid any ambiguity or conflict. Thus usability is a crucial issue. Especially as we cannot assume that every data provider would be an expert in privacy or any type of technology, GUI tools that are intuitive and instructive must be provided for them. We believe that many valuable lessons can be learned from existing technology related to P3P [20] and APPEL [19] and work on user interaction design (See [21] for example). It is important that data providers have a clear understanding of the guarantees provided by each privacy level.

Data generalization. Needless to say, devising a quality data generalization technique is one of the key challenges. There are two important issues to be addressed here. The first issue is that the generalization process must preserve meaningful information from actual data as inadequate information would not be of any use. For example, although numeric or structured data may be easy to generalize into multiple levels that are meaningful, it is unclear how unstructured data (e.g., textual data) should be generalized into multiple levels. We need also to devise generalization policies and ontologies supporting systematic and consistent data generalization across the database. The other important issue is that the generalization process must produce a sufficient level of data privacy by effectively suppressing distinctive information in individual data. For instance, consider the names of individuals. There are certain names that are less frequent than the others, and inadequate generalization techniques would not hide the uniqueness of such names. Moreover, if the content of database may changes dynamically, hiding such uniqueness becomes much more challenging. Clearly, a key challenge in data generalization is to balance the trade-off between information preservation and information loss. Also, generalization must be efficient. In many cases, the system will have to perform data generalization “on the fly” while processing a query. In other cases, post-processing of queries is required because what has to be returned may depend on various factors such as the cardinality and statistics of the results or even the past accesses. Many valuable lessons can be learned from var-

ious generalization techniques that are available in statistical databases [1]. The main challenge here is that these techniques may have to be used dynamically in a variety of settings, ranging from data storage to query processing.

Metrics for data privacy and data usability. So far, we have claimed that both privacy and usability of data can be achieved when data is sufficiently generalized. However, a key question is: how can we determine whether or not a certain generalization strategy provides a sufficient level of privacy and usability? As one can generalize data in various ways and degrees, we need metrics that methodologically measure the privacy and usability of generalized data. Such metrics are necessary to devise generalization techniques that satisfy the requirements of both data providers and data users.

Metadata storage. In our “naive” model, we assumed that the collected data is generalized and stored into multiple privacy levels at the preprocessing stage. This approach is simple and effective, yet may require large storage space. For instance, suppose there are n privacy levels in a system. This means that the required storage space would be n times the size of the collected data. Another approach is to postpone the generalization process to the time of data access. This method does not require any additional storage and may help avoid unnecessary data generalization. As some data items are accessed much less frequently than the others, those rarely accessed data items do not have to be generalized unless they are accessed. However, the overall performance may significantly suffer. Another possible solution is to use both pre-generalization and post-generalization selectively. For example, only data items that are expected to be accessed frequently are pre-generalized and stored. Other data items that are not pre-processed should be generalized when they are accessed. Also, for better performance the post-generalized data may be cached in a special space. Using this approach, one can try to reduce the overall cost of generalization process. However, a challenge here is to balance the trade-off between storage and performance. Yet another approach could be based on the use of views, which would have to be extended with innovative capabilities for value generation and transformation.

Complex query processing. In this paper we have considered only simple queries; i.e., queries without join, sub-queries or aggregations. A key question here is whether complex queries can be introduced in our model. Even though it seems that they can be correctly processed in the model, it is not clear whether the results of such queries would be still meaningful. For instance, how do we calculate the sum of several generalized values, and how do we interpret such results?

Applicability to general-purpose access control. Although we have limited our discussion to access control for privacy protection, we believe it is possible to extend our model to a general-purpose access control model. For instance, each user can be assigned a trust level⁵, and the access control system can control, based on the user’s trust level, degrees of precision on accessible information. This

⁵The trust level is not chosen by users, but assigned to users by security officers.

approach is very similar to multilevel secure database systems [6, 15, 13], where every piece of information is classified into a security level and every user is assigned a security clearance. However, the main difference is that our approach can provide a much finer level of control as the access control decision is based on the question of “how much information can be allowed for a certain user”, rather than “is information allowed for a certain user or not”. In other words, our model utilizes the elaborated version of cover story in multilevel secure databases. This type of finer grained access control can be extremely useful for internal access control within an organization as well as information sharing between organizations. Even though such an extension seems very promising at this point, further investigation is required to confirm this hypothesis.

Other issues. There are many other issues that require careful investigation, such as problems of polyinstantiation [11, 14], inference and integrity. Addressing such issues is also crucial for the development of comprehensive access control models for high-assurance privacy.

4. RELATED WORK

To date, several approaches have been reported that deal with various aspects of the problem of high-assurance privacy systems. Here we briefly discuss the approaches that have provided some initial solutions that can certainly be generalized and integrated into comprehensive solutions to such problem.

The W3C’s Platform for Privacy Preference (P3P) [20] allows web sites to encode their privacy practice, such as what information is collected, who can access the data for what purposes, and how long the data will be stored by the sites, in a machine-readable format. P3P enabled browsers can read this privacy policy automatically and compare it to the consumer’s set of privacy preferences which are specified in a privacy preference language such as A P3P Preference Exchange Language (APPEL) [19], also designed by the W3C.

The concept of Hippocratic databases that incorporates privacy protection within relational database systems was introduced by Agrawal et al. [2]. The proposed architecture uses privacy metadata, which consist of privacy policies and privacy authorizations stored in two tables. A privacy policy defines for each attribute of a table the usage purpose(s), the external-recipients and a retention period, while a privacy authorization defines which purposes each user is authorized to use.

Byun et al. presented a comprehensive approach for privacy preserving access control based on the notion of purpose [4, 5]. In the model, purpose information associated with a given data element specifies the intended use of the data element, and the model allows multiple purposes to be associated with each data element. The granularity of data labeling is discussed in detail in [4], and a systematic approach to implement the notion of access purposes, using roles and role-attributes is presented in [5].

Previous work on multilevel secure relational databases [6, 13, 15] also provides many valuable insights for designing a fine-grained secure data model. In a multilevel relational database system, every piece of information is classified into a security level, and every user is assigned a security clear-

ance. Based on this access class, the system ensures that each user gains access to only the data for which he has proper clearance, according to the basic restrictions. These constraints ensure that there is no information flow from a lower security level to a higher security level and that subjects with different clearances see different versions of multilevel relations.

In order to prevent re-identification of anonymized data, Sweeney introduced the notion of k-anonymity [18]. K-anonymity requires that information about each individual in a data release be indistinguishable from at least k-1 other individuals with respect to a particular set of attributes. Sweeney also proposed a technique using generalization and suppression of data to achieve k-anonymity with minimal distortion [17].

5. CONCLUSIONS

In this paper, we discussed a new approach for access control that maximizes the usability of private information for enterprises while, at the same time, assuring privacy. We believe that one direction for next-generation DBMS technology is represented by DBMS with high-assurance security and privacy. The “naive” model we presented in this paper provides an example of access control for such a new DBMS. Based on this model, we discussed many challenges that need to be addressed. We conclude this paper by saying that ultimately suitable access control systems with high-privacy assurance will be built by integrating techniques such as view mechanisms, statistical databases, anonymization, privacy-preserving computation and data mining. The main challenge is how to integrate such techniques in a full-fledged DBMS ensuring good performance.

6. REFERENCES

- [1] Nabil Adam and John Wortmann. Security-control methods for statistical databases: A comparative study. *ACM Computing Surveys (CSUR)*, 21, 1989.
- [2] Rakesh Agrawal, Jerry Kiernan, Ramakrishnan Srikant, and Yirong Xu. Hippocratic databases. In *The 28th International Conference on Very Large Databases (VLDB)*, 2002.
- [3] Paul Ashley, Calvin S. Powers, and Matthias Schunter. Privacy promises, access control, and privacy management. In *Third International Symposium on Electronic Commerce*, 2002.
- [4] Jiwon Byun, Elisa Bertino, and Ninghui Li. Purpose based access control for privacy protection in relational database systems. Technical Report 2004-52, Purdue University, 2004.
- [5] Jiwon Byun, Elisa Bertino, and Ninghui Li. Purpose based access control of complex data for privacy protection. In *Symposium on Access Control Model And Technologies (SACMAT)*, 2005.
- [6] Dorothy Denning, Teresa Lunt, Roger Schell, William Shockley, and Mark Heckman. The seaview security model. In *The IEEE Symposium on Research in Security and Privacy*, 1998.
- [7] Xin Dong, Alon Halevy, Jayant Madhavan, and Ema Nemes. Reference reconciliation in complex information spaces. In *ACM International Conference on Management of Data (SIGMOD)*, 2005.
- [8] I. P. Fellegi and A. B. Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 1969.
- [9] IBM. *The Enterprise Privacy Authorization Language (EPAL)*. Available at www.zurich.ibm.com/security/enterprise-privacy/epal.
- [10] Kristen LeFevre, Rakesh Agrawal, Vuk Ercegovic, Raghu Ramakrishnan, Yirong Xu, and David DeWitt. Disclosure in hippocratic databases. In *The 30th International Conference on Very Large Databases (VLDB)*, August 2004.
- [11] Fausto Rabitti, Elisa Bertino, Won Kim, and Darrell Woelk. A model of authorization for next-generation database systems. In *ACM Transactions on Database Systems (TODS)*, March 1991.
- [12] Forrester Research. Privacy concerns cost e-commerce \$15 billion. Technical report, September 2001. Available at www.forrester.com.
- [13] Ravi Sandhu and Fang Chen. The multilevel relational data model. In *ACM Transactions on Information and System Security*, 1998.
- [14] Ravi Sandhu and Sushil Jajodia. Polyinstantiation integrity in multilevel relations. In *IEEE Symposium on Security and Privacy*, 1990.
- [15] Ravi Sandhu and Sushil Jajodia. Toward a multilevel secure relational data model. In *ACM International Conference on Management of Data (SIGMOD)*, 1991.
- [16] S. Sarawagi and A. Bhamidipaty. Interactive deduplication using active learning. In *ACM International conference on Knowledge discovery and data mining (SIGKDD)*, 2002.
- [17] Latanya Sweeney. Achieving k-anonymity privacy protection using generalization and suppression. In *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 2002.
- [18] Latanya Sweeney. K-anonymity: A model for protecting privacy. In *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 2002.
- [19] World Wide Web Consortium (W3C). *A P3P Preference Exchange Language 1.0 (APPEL 1.0)*. Available at www.w3.org/TR/P3P-preferences.
- [20] World Wide Web Consortium (W3C). *Platform for Privacy Preferences (P3P)*. Available at www.w3.org/P3P.
- [21] Kaping Yee. User interaction design for secure systems. In *The 4th International Conference on Information and Communications Security*, 2002.