

Paper ID: CSEIT04

## DATA LINEAGE IN MALICIOUS ENVIRONMENTS

Mrs. Nilima Dheeraj Patil,  
Prof. N.P.Karlekar

Department of Computer Engineering, Sinhgad Institute of Technology, Lonavala, Pune

**Abstract:** In this paper cover Data Lineage in Malicious Environment. A data provider has given precise data to a set of supposedly trusted node protocol. Some of the data are leaked and found in an unjustified place. The provider must assess the nearest that the crevice data came from one or more node protocol, as opposed to having been individually fetch by other terms. We propose data allocation strategies that improve the probability of identifying crevices. These methods do not build on alterations of the released data. In some cases, we can also implant “realistic but wrong and threaten data” data records to further improve our chances of detecting crevice and identifying the guilty party. While sending data over the network there is lots of illegitimate user trying to get useful information. There should be proper security should be provided to data which is send to network. Now a day’s new tech cell phones use have been increased rapidly and the applications used in new tech cell phones can get easy access to our confidential information. So for avoiding this we used the data lineage mechanism. We give the wrong and threaten data information to guilty agent. We design and analyze a new accountable data transfer protocol between two entities within a malicious environment by building upon oblivious transfer, robust Watermarking, and signature primitives. Finally, we execute an experimental evaluation to demonstrate the practicality of our protocol and apply our framework to the vital data leakage case of data outsourcing and social networks. In general, we consider our lineage framework for data transfer, to be a key step towards achieving accountability by design.

**Keyword:** Data Leakage Prevention, Data Privacy Leakage Model Watermarking, Data Leakage Protection, Data Loss Prevention.

### I. INTRODUCTION

Data Leakage is a vital concern for the establish farms in this increasingly networked world these days. Illegitimate disclosure may have serious consequences for an farm in both long term and short term. Risks include losing clients and stakeholder Confidence, tarnishing of brand image, landing in undesirable lawsuits, and overall losing goodwill and market share in the, industry. To prevent from all these unwanted and nasty activities from happening, an organized effort is needed to control and monitor the information flow inside and outside the farm.

Here is our attempt to demystify the jargon surrounding the data leakage prevention procedures which will help you to choose and apply the best suitable option for your own business. Leakage describes an unwanted loss of something which escapes from its proper location and Lineage describes as data flow across multiple entities that take two characteristic, principal roles (i.e., owner and consumer). We define the exact security guarantees required by such a data lineage mechanism toward identification of a guilty entity, and identify the simplifying non-repudiation and honesty assumptions. In the course of doing business, sometimes sensitive data must be handed over to supposedly trusted third parties. For example, a hospital may give patient records to researchers who will devise new treatments. Similarly, a company may have partnerships with other companies that require sharing customer data. Another enterprise may outsource its data processing, so data must be given to various other companies. The owner of the data can be called as provider and the supposedly trusted third parties the node protocol. The goal is to detect when the providers sensitive data have been leaked by node protocol, and if possible to identify the agent that crevice the data.

### II. OVERVIEW OF DATA LINEAGE

Data Leakage Prevention is the category of solutions which help a farm to apply control and monitors for preventing the unwanted accidental or malicious leakage of precise information to illegitimate entities in or outside the farm. Here sensitive information may refer to farm’s internal process documents, strategic business plans, intellectual property, financial statements, security policies, network diagrams, blueprints etc.

#### 2.1. Need Data Lineage

There are many areas where data leakage may occur, so it is very essential detect such kind of detection, following users may lead to data leakage-

##### 1. The security illiterate

- Majority of employees with little or no knowledge of security

- Corporate risk because of accidental breaches

##### 2. The gadget needs

- Introduce a variety of devices to their work PCs

- Download software

##### 3. The unlawful residents

- Use the company IT resources in ways they shouldn’t

- i.e. by storing music, movies, or playing games

##### 4. The malicious/disgruntled employees

K.E. Society's

**RAJARAMBAPU INSTITUTE OF TECHNOLOGY**

- Typically minority of employees
- Gain access to areas of the IT system to which they shouldn't Send corporate data (e.g., customer lists, RD, etc.) to third parties.

## 2.2. Generic Data Leakage Prevention

### 1 Deploy Security Mechanisms

- a). Firewalls, IDS's & antivirus software
- b).Thin-client architecture

### 2 Advanced Security Measures

- a). Use of pattern based monitoring tools
- b).Use of reasoning algorithms

### 3 Access Control and monitor & Encryption

- a).Access Control and monitor & Device Control and monitor
- b).Storage of encryption keys

## III. RELATED WORK- CREATING ENCRYPTED DIGITAL WATERMARK

Our process and watermarking are similar in the sense of providing node protocol with some kind of receiver identifying information. However, by its very nature, a watermark modifies the item being watermarked. If the object to be watermarked cannot be modified, then a watermark cannot be inserted. In such cases, methods that attach watermarks to the distributed data are not applicable. Finally, there are also lots of other works on mechanisms that allow only authorized users to access sensitive data through access control and monitor policies. Such processes prevent in some sense dataleakage by sharing information only with trusted parties. However, these policies are restrictive and may make it impossible to satisfy node protocol requests. LIME (Lineage In the Malicious Environment) can be used with any type of data for which watermarking schemes exist. Therefore, we briefly describe different watermarking techniques for different data types. Most watermarking schemes are designed for multimedia files such as images, videos, and audio files. In these multimedia files, watermarks are usually embedded by using a transformed representation (e.g. discrete cosine, wavelet or Fourier transform) and modifying transform domain coefficients. Watermarking techniques have also been designed for other data types such as relational databases, text files and even Android apps. The first two are especially interesting, as they allow us to apply LIME to user databases or medical records. Watermarking relational databases can be done in different ways. The most common solutions are to embed information in noise-tolerant attributes of the entries or to create wrong and threaten data database entries. For watermarking of texts, there are two main processes. The first one embeds information by changing the text's appearance (e.g. changing distance between words and lines) in a way that is imperceptible to humans. The second process is also referred to as language watermarking and

works on the semantic level of the text rather than on its appearance. A mechanism also has been proposed to insert watermarks to Android apps. This mechanism encodes a watermark in a permutation graph and hides the graph as a linked list in the application. Due to the list representation, watermarks are encoded in the execution state of the application rather than in its syntax, which makes it robust against attacks. In this process the authors propose to rather remove existing information than adding new information or modifying existing information. Thereby the watermarking scheme guarantees that no false entries are introduced. The above schemes can be employed in our framework to create data lineage for documents of the respective formats. The only Modification that might be necessary when applying our scheme to a different document type is the splitting algorithm. For example for images it makes more sense to take small rectangles of the original image instead of simply taking the consecutive bytes from the pixel array. Embedding multiple watermarks into a single document has been discussed in literature and there are different techniques available. In they discuss multiple rewatermarking and in the focus is on segmented watermarking. Both projects show in experimental results that multiple watermarking is possible which is very vital for our scheme, as it allows us to create a lineage over multiple levels. It would be desirable not to reveal the private watermarking key to the auditor during the auditor's investigation, so that it can be safely reused, but as discussed in current public key watermarking schemes are not secure and it is doubtful if it is possible to design one that is secure. In Sadeghi presents processes to zero knowledge watermark detection. With this technology it is possible to convince another party of the presence of a watermark in a document without giving any information about the detection key or the watermark itself. However, the scheme discussed in also hides the content of the watermark itself and is therefore unfit for our case, as the auditor has to know the watermark to identify the guilty person. Furthermore, using a technology like this would come with additional constraints for the chosen watermarking scheme.

## IV.APPLICATION

It involves study of unobtrusive techniques for detecting leakage of a set of objects or records. Specifically, following scenario can be studied: After giving a set of objects to node protocol, the provider discovers some of those same objects in an unauthorized place. At this point, the provider can assess the nearest that the leaked data came from one or more node protocol, as opposed to having been independently fetch by other terms. In the proposed process, a model is designed for assessing the guilt of node protocol. The algorithms are also presented for distributing objects to node protocol, in a way that improves the chances of identifying a leaker.

Finally, the option of adding wrong and threaten data objects to the distributed set is also considered. Such objects do not correspond to real entities but appear realistic to the node protocol. In a sense, the wrong and threaten data objects act as a type of watermark for the entire set, without modifying any individual members. If it turns out that an agent was given one or more wrong and threaten data objects that were leaked, then the provider can be more confident that agent was guilty. In the Proposed System, the hackers can be traced with good amount of evidence.

#### V. CONCLUSION & FUTURE SCOPE

We present LIME, a model for countable data transfer across multiple entities. We define participating parties, their interrelationships and give a concrete instantiation for a data transfer protocol using a new combination of oblivious transfer, robust watermarking and digital signatures. Although LIME does not actively prevent data leakage, it introduces reactive accountability. Thus, it will deter malicious parties from leaking private documents and will encourage honest (but careless) parties to provide the required protection for sensitive data. LIME is flexible as we differentiate between trusted senders (usually owners) and untrusted senders (usually consumers). In the case of the trusted sender, a very simple protocol with little overhead is possible. The untreated sender requires a more complicated protocol, but the results are not based on trust assumptions and therefore they should be able to convince a neutral entity (e.g. a judge). Our work also motivates further research on data leakage detection techniques for various document types and case. For example, it will be an interesting future research direction to design a verifiable lineage protocol for derived data.

#### ACKNOWLEDGEMENT

For all the efforts behind the paper work, I first & foremost would like to express my sincere appreciation to the staff of Dept. of Computer Engineering, Sinhgad Institute of Technology, Lonavala, Pune. For their extended help & suggestions at every stage of this paper. It is with a great sense of gratitude that I acknowledge the support, time to time suggestions and highly indebted to my guide Prof. N. P. Karlekar. Finally, I pay sincere thanks to all those who indirectly and directly helped me towards the successful completion of the paper.

#### REFERENCES

[1] "Chronology of data breaches," <http://www.privacyrights.org/data-breach>.  
[2] "Data breach cost," <http://www.symantec.com/about/news/release/article.jsp?prid=2011030801>.

[3] "Privacy rights clearinghouse," <http://www.privacyrights.org>.  
[4] "Electronic Privacy Information Center (EPIC)," <http://epic.org>, 1994.  
[5] "Facebook in Privacy Breach," <http://online.wsj.com/article/SB10001424052702304772804575558484075236968.html>.  
[6] "Offshore outsourcing," [http://www.computerworld.com/s/article/19938/Offshore\\_outsourcing\\_cited\\_in\\_Florida\\_data\\_leak](http://www.computerworld.com/s/article/19938/Offshore_outsourcing_cited_in_Florida_data_leak).  
[7] A. Mascher-Kampfer, H. Stogner, and A. Uhl, "Multiplere-watermarking case," in Proceedings of the 13<sup>th</sup> International Conference on Systems, Signals, and Image Processing (IWSSIP 2006).Citeseer, 2006, pp. 53–56.  
[8] P. Papadimitriou and H. Garcia- Molina, "Data leakage detection," Knowledge and Data Engineering, IEEE Transactions on, vol. 23, no. 1, pp. 51–63, 2011.  
[9] "Pairing-Based Cryptography Library (PBC)," <http://crypto.stanford.edu/pbc>.  
[10] I. J. Cox, J. Kilian, F. T. Leighton, and T. Shamoon, "Secure spread spectrum watermarking for multimedia," Image Processing, IEEE Transactions on, vol. 6, no. 12, pp. 1673–1687, 1997.  
[11] B. Pfitzmann and M. Waidner, "Asymmetric fingerprinting for larger collusions," in Proceedings of the 4<sup>th</sup> ACM conference on Computer and communications security, ser. CCS '97, 1997, pp. 151–160.  
[12] S. Goldwasser, S. Micali, and R. L. Rivest, "A digital signature scheme secure against adaptive chosen-message attacks," SIAM J. Comput., vol. 17, no. 2, pp. 281–308, 1988.  
[13] A. Adelsbach, S. Katzenbeisser, and A.-R. Sadeghi, "A computational model for watermark robustness," in Information Hiding Springer, 2007, pp. 145–160.  
[14] J. Kilian, F. T. Leighton, L. R. Matheson, T. G. Shamoon, R. E. Tarjan, and F. Zane, "Resistance of digital watermarks to collusive attacks," in IEEE International Symposium on Information Theory, 1998, pp. 271–271.  
[15] M. Naor and B. Pinkas, "Efficient oblivious transfer protocols," in Proceedings of the Twelfth Annual ACM Symposium on Discrete Algorithms, 2001, pp. 448–457.  
[16] "GNU Multiple Precision Arithmetic Library (GMP)," <http://gmplib.org/>.  
[17] D. Boneh, B. Lynn, and H. Shacham, "Short signatures from the Weil pairing," in Advances in Cryptology-ASIACRYPT2001. Springer, 2001, pp. 514–532.  
[18] W. Dai, "Crypto++ Library," <http://cryptopp.com>.  
[19] P. Meerwald, "Watermarking toolbox," [http://www.cosy.sbg.ac.at/\\_pmeerw/Watermarking/source](http://www.cosy.sbg.ac.at/_pmeerw/Watermarking/source).