



RESEARCH ARTICLE

DATA LEAKAGE DETECTION

Ms. N. Bangar Anjali¹, Ms. P. Rokade Geetanjali², Ms. Patil Shivilila³, Ms. R. Shetkar Swati⁴, Prof. N B Kadu⁵

¹BE Computer, Pravara Rural Engineering College, Loni, Tal: Rahata, Dist: A, Nagar, Pin: 41373, India

²BE Computer, Pravara Rural Engineering College, Loni, Tal: Rahata, Dist: A, Nagar, Pin: 41373, India

³BE Computer, Pravara Rural Engineering College, Loni, Tal: Rahata, Dist: A, Nagar, Pin: 41373, India

⁴BE Computer, Pravara Rural Engineering College, Loni, Tal: Rahata, Dist: A, Nagar, Pin: 41373, India

⁵M.Tech Computer, Pravara Rural Engineering College, Loni, Tal: Rahata, Dist: A, Nagar, Pin: 41373, India

¹ goodluckanju@gmail.com, ² rokadegitanjali@gmail.com, ³ Shivilila21@gmail.com,

⁴ Swatishetkar51@gmail.com, ⁵ kamleshkadu@rediffmail.com

Abstract— *This paper contains the implementation of data leakage detection model's technology. This paper deals with new technique of research for secured data transmission & leakage detection, if it gets leaked. A data distributor has given sensitive data to the trusted third parties (agents).if some data are leaked and found in an unauthorized place. The distributor must analyze that the leaked data came from (where) one or more agents. We propose data allocation strategies to improve the probability of identifying leakages. We can also add "fake object" to further improve our chances of detecting leakage and identifying the guilty party.*

Key Terms: - *data leakage; implicit request; explicit request; steganography; guilty agent; AES algorithm*

I. INTRODUCTION

Data leakage is the unauthorized transmission of data or information from within an organization to an external destination or recipient. Data leakage is defined as the accidental or intentional distribution of private or sensitive data to an unauthorized entity. Sensitive data of companies and organization includes intellectual property, financial information, patient information, personal credit card data and other information depending upon the business and the industry. A data distributor has given this sensitive data to a set of supposedly trusted agents (third parties). Some of the data are leaked and found in an unauthorized place. The distributor must assess the likelihood that the leaked data came from one or more agents, as opposed to having been independently gathered by other means. We call the owner of the data the distributor and the supposedly trusted third parties the agents. Our goal is to detect when the distributor's sensitive data have been leaked by agents, and if possible to identify the agent that leaked the data. We propose data allocation strategies that improve the probability of identifying leakages. These methods do not rely on alterations of the released data (e.g., watermarks). In some cases, we can also inject —realistic but fake|| data records to further improve our chances of detecting leakage and identifying the guilty party.

II. PROBLEM SETUP AND NOTATION

A. Entities and Agents

Let the distributor database owns a set $S = \{t_1, t_2, \dots, t_m\}$ which consists of data objects. Let the no of agents be A_1, A_2, \dots, A_n . The distributor distributes a set of records S to any agents based on their request such as sample or explicit request.

- Sample request $R_i = \text{SAMPLE}(T, m_i)$: Any subset of m_i records from T can be given to U_i [1].
- Explicit request $R_i = \text{EXPLICIT}(T; \text{condi})$: Agent U_i receives all T objects that satisfy condition.

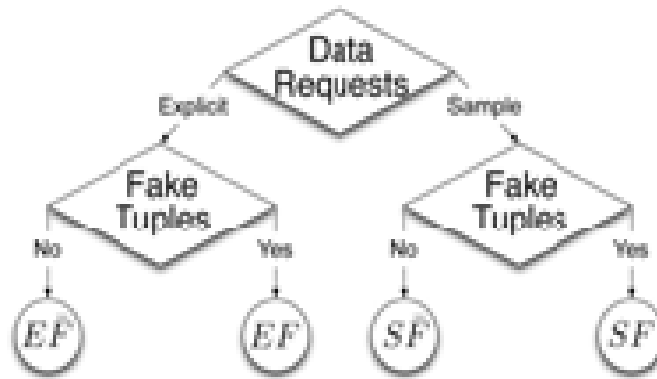


Fig.1 Leakage Problem Instances

The objects in T could be of any type and size, e.g. they could be tuples in a relation, or relations in a database. After giving objects to agents, the distributor discovers that a set S of T has leaked. This means that some third party called the target has been caught in possession of S . For example, this target may be displaying S on its web site, or perhaps as part of a legal discovery process, the target turned over S to the distributor. Since the agents (A_1, A_2, \dots, A_n) have some of the data, it is reasonable to suspect them leaking the data. However, the agents can argue that they are innocent, and that the S data was obtained by the target through other means.

B. Guilty Agents

Guilty agents are the agents who had leaked the data. Suppose the agent say A_i had leaked the data knowingly or unknowingly. Then automatically notification will be the send to the distributor defining that agent A_i had leaked the particular set of records which also specifies sensitive or non-sensitive records. Our goal is to estimate the likelihood that the leaked data came from the agents as opposed to other sources.

III. RELATED WORK

The guilt detection approach we present is related to the data provenance problem [1] : tracing the lineage of an S object implies essentially the detection of the guilty agents. [2]It provides a good overview on the research conducted in this field. Suggested solutions are domain specific, such as lineage tracing for data Warehouses [3], and assume some prior knowledge on the way a data view is created out of data sources. Our problem formulation with objects and sets is more general and simplifies lineage tracing, since we do not consider any data transformation from R_i sets to S .As far as the data allocation strategies are concerned, our work is mostly relevant to watermarking that is used as a means of establishing original ownership of distributed objects. Watermarks were initially used in images [4], video [5] and audio data [6] whose digital representation includes considerable redundancy. Our approach and watermarking are similar in the sense of providing agents 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 policies[7],[8]. Such approaches prevent in some sense data leakage by sharing information only with trusted parties. However, these policies are restrictive and may make it impossible to satisfy agents' requests.

IV. RESULTS OF DATA LEAKAGE DETECTION MODEL

A. Agent Guilt Model

To compute this $PrfGijSg$, we need an estimate for the probability that values in S can be "guessed" by the target. For instance, say some of the objects in T are emails of individuals. We can conduct an experiment and ask a person with approximately the expertise and resources of the target to find the email of say 100 individuals. If this person can find say 90 emails, then we can reasonably guess that the probability of finding one email is 0.9. On the other hand, if the objects in question are bank account numbers, the person may only discover say 20, leading to an estimate of 0.2. We call this estimate p_t , the probability that object t can be guessed by the target. To simplify the formulas that we present in the rest of the paper, we assume that all T objects have the same p_t , which we call p . Our equations can be easily generalized to diverse p_t 's though they become cumbersome to

display. Next, we make two assumptions regarding the relationship among the various leakage events. The first assumption simply states that an agent's decision to leak an object is not related to other objects.

B. Guilt Model Analysis

In order to see how our model parameters interact and to check if the interactions match our intuition, in this section, we study two simple scenarios. In each scenario, we have a target that has obtained all the distributor's objects, i.e., $T \setminus S$.

B.1 Impact of Overlap between R_i and S

In this section, we again study two agents, one receiving all the $T \setminus S$ data and the second one receiving a varying fraction of the data. The probability of guilt for both agents, as a function of the fraction of the objects owned by U_2 , i.e., as a function of $j \in R_2 \setminus S_j = j \setminus S_j$. In this case, p has a low value of 0.2, and U_1 continues to have all $16S$ objects. Note that in our previous scenario, U_2 has 50 percent of the S objects. We see that when objects are rare ($p \setminus 0.2$), it does not take many leaked objects before we can say that U_2 is guilty with high confidence. This result matches our intuition: an agent that owns even a small number of incriminating objects is clearly suspicious. The same scenario, except for values of p equal to 0.5 and 0.9. We see clearly that the rate of increase of the guilt probability decreases as p increases. This observation again matches our intuition: As the objects become easier to guess, it takes more and more evidence of leakage (more leaked objects owned by U_2) before we can have high confidence that U_2 is guilty. In, we study an additional scenario that shows how the sharing of S objects by agents affects the probabilities that they are guilty. The scenario conclusion matches our intuition: with more agents holding the replicated leaked data, it is harder to lay the blame on any one agent.

V. EXISTING SYSTEM

We consider applications where the original sensitive data cannot be perturbed. Perturbation is a very useful technique where the data are modified and made "less sensitive" before being handed to agent. In some cases it is important not to alter the original distributor's data. Traditionally, leakage detection is handled by watermarking, e.g., a unique code is embedded in each distributed copy. If that copy is later discovered in the hands of an unauthorized party, the leaker can be identified. Watermarks can be very useful in some cases, but again, involve some modification of the original data. Furthermore, watermarks can sometimes be destroyed if the data recipient is malicious. *E.g.* 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. We call the owner of the data the distributor and the supposedly trusted third parties the agents.

VI. PROPOSED SYSTEM

Our goal is to detect when the distributor's sensitive data has been leaked by agents, and if possible to identify the agent that leaked the data. We develop *unobtrusive* techniques for detecting leakage of a set of objects or records. In this section we develop a model for assessing the "guilt" of agents. We also present algorithms for distributing objects to agents, in a way that improves our chances of identifying a leaker. Finally, we also consider the option of adding "fake" objects to the distributed set. Such objects do not correspond to real entities but appear realistic to the agents. In a sense, the fake objects acts as a type of watermark for the entire set, without modifying any individual members. If it turns out an agent was given one or more fake objects that were leaked, then the distributor can be more confident that agent was guilty.

6.1 Watermarking:

We describe a digital watermarking method for use in audio, image, video and multimedia data.

The watermark is difficult for an attacker to remove, even when several individuals conspire together with independently watermarked copies of the data. It is distortions such as digital-to-analog and analog-to-digital conversion, resampling, quantization, dithering, compression, rotation, translation, cropping and scaling. The same digital watermarking algorithm can be applied to all three media under consideration with only minor modifications, making it especially appropriate for multimedia products.

6.2 Steganography

Steganography is a technique for hiding a secret message within a larger one in such a way that others can't discern the presence or contents of the hidden message. A plain text message may be hidden in one of two ways, the method of steganography conceal the existence of the message, whereas the outsiders of cryptography render the message unintelligible to transformation of the text. Steganography serves as a means for private, secure and sometimes malicious communication.

VII. MODULES OF DATA LEAKAGE DETECTION SYSTEM

A. Data Allocation Module

The main focus of our project is the data allocation problem as how can the distributor intelligently give data to agents in order to improve the chances of detecting a guilty agent, Admin can send the files to the authenticated user, users can edit their account details etc. Agent views the secret key details through mail. In order to increase the chances of detecting agents that leak data.

B. Fake Object Module

The distributor creates and adds fake objects to the data that he distributes to agents. Fake objects are objects generated by the distributor in order to increase the chances of detecting agents that leak data. The distributor may be able to add fake objects to the distributed data in order to improve his effectiveness in detecting guilty agents.

C. Optimization Module

The Optimization Module is the distributor's data allocation to agents has one constraint and one objective. The agent's constraint is to satisfy distributor's requests, by providing them with the number of objects they request or with all available objects that satisfy their conditions. His objective is to be able to detect an agent who leaks any portion of his data. User can able to lock and unlock the files for secure.

D. Data Distributor Module

A data distributor has given sensitive data to a set of supposedly trusted agents (third parties). Some of the data is leaked and found in an unauthorized place. The distributor must assess the likelihood that the leaked data came from one or more agents, as opposed to having been independently gathered by other means Admin can able to view the which file is leaking and fake user's details also.

E. Agent Guilt Module

Probability of guilt $Pr\{Gi|S\}$ can be computed by estimating the probability that the target can guess objects in "S". The proposed guilt model makes two assumptions. The first assumption is that the source of a leaked object can be of any agent. The second assumption is that an object which is part of set of objects distributed can only be obtained from one of the agents or through other means. With these assumptions the probability of guilt is computed as $Pr\{Ui\text{ leaked t to } S\} = \{ 1-p, \text{ if } Ui \in Vt$
 $|Vt|$

0, otherwise

VIII. DATA ALLOCATION STRATEGIES

The data allocation strategies used to solve the problem of data distribution as discussed in previous sections exactly or approximately are provided in the form of various algorithms. The algorithms are provided here.

8.1 Explicit Data Request

Algorithm 1 Allocation for Explicit Data Requests (EF)

```

Input:  $R_1, \dots, R_n, cond_1, \dots, cond_n, b_1, \dots, b_n, B$ 
Output:  $R_1, \dots, R_n, F_1, \dots, F_n$ 
1:  $R \leftarrow \phi$  ▷ Agents that can receive fake objects
2: for  $i=1, \dots, n$  do
3: if  $b_i > 0$  then
4:  $R \leftarrow R \cup \{i\}$ 
5:  $F_i \leftarrow \phi$ 
6: while  $B > 0$  do
7:  $i \leftarrow \text{SELECTAGENT}(R, R_1, \dots, R_n)$ 
8:  $f \leftarrow \text{CREATEFAKEOBJECT}(R_i, F_i, cond_i)$ 
9:  $R_i \leftarrow R_i \cup \{f\}$ 
10:  $F_i \leftarrow F_i \cup \{f\}$ 
11:  $b_i \leftarrow b_i - 1$ 
12: if  $b_i = 0$  then
13:  $R \leftarrow R \setminus \{R_i\}$ 
14:  $B \leftarrow B - 1$ 
    
```

Fig.2 – Allocation for explicit data requests

It is a general algorithm that is used by other algorithms.

8.2 Sample Data Request

This algorithm is meant for making a greedy choice of choosing an agent that causes improvement in the sum-objective.

Algorithm 4: Allocation for Sample Data Requests(SF)

```

Input:  $m_1, \dots, m_n, |T|$           ▷ Assuming  $m_i \leq |T|$ 
Output:  $R_1, \dots, R_n$ 
1:  $a \leftarrow 0_{|T|}$           ▷  $a[k]$ : number of agents who have received object  $t_k$ 
2:  $R_1 \leftarrow \phi, \dots, R_n \leftarrow \phi$ 
3: remaining  $\leftarrow \sum_{i=1}^n m_i$ 
4: while remaining  $> 0$  do
5: for all  $i=1, \dots, n: |R_i| < m_i$  do
6:  $k \leftarrow \text{SELECTOBJECT}(i, R_i)$  ▷ May also use additional parameters
7:  $R_i \leftarrow R_i \cup \{t_k\}$ 
8:  $a[k] \leftarrow a[k] + 1$ 
9: remaining  $\leftarrow \text{remaining} - 1$ 
    
```

Fig. 3– Allocation for sample data requests

8.3 AES Algorithm

The AES algorithm is based on permutations and substitutions. Permutations are rearrangements of data, and substitutions replace one unit of data with another.

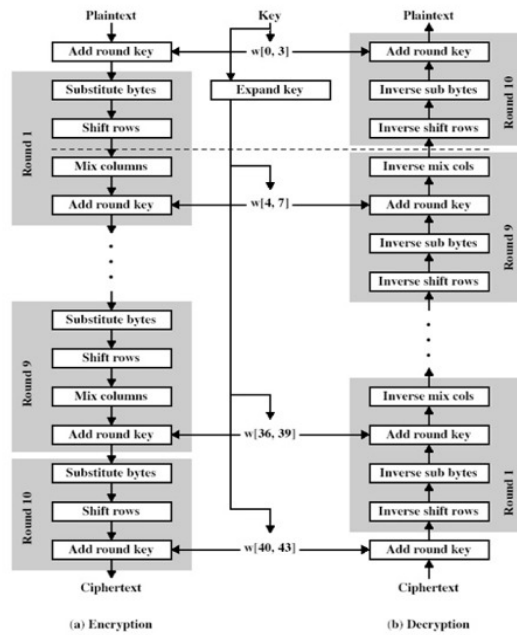


Fig.4 The Advanced Encryption Algorithm

AES performs permutations and substitutions using several different techniques. The four operations SubBytes, Shift SiftRows, MixColumns, and AddRoundKey are called inside a loop that executes N_r times—the number of rounds for a given key size, less 1. The number of rounds that the encryption algorithm uses is either 10, 12, or 14 and depends on whether the seed key size is 128, 192, or 256 bits. In this example, because N_r equals 12, the four operations are called 11 times. After this iteration completes, the encryption algorithm finishes by calling SubBytes, ShiftRows, and AddRoundKey before copying the State matrix to the output parameter. In summary, there are four operations that are at the heart of the AES encryption algorithm. AddRoundKey substitutes groups of 4 bytes using round keys generated from the seed key value. SubBytes substitutes individual bytes using a substitution table. ShiftRows permutes groups of 4 bytes by rotating 4-byte rows. Mix Columns substitutes bytes using a combination of both field addition and multiplication.

IX. CONCLUSION

From this study we conclude that the data leakage detection system model is very useful as compare to the existing watermarking model. We can provide security to our data during its distribution or transmission and even we can detect if that gets leaked. Thus, using this model security as well as tracking system is developed. Watermarking can just provide security using various algorithms through encryption, whereas this model provides security plus detection technique.

Our model is relatively simple, but we believe that it captures the essential trade-offs. The algorithms we have presented implement a variety of data distribution strategies that can improve the distributor's chances of identifying a leaker. We have shown that distributing objects judiciously can make a significant difference in identifying guilty agents, especially in cases where there is large overlap in the data that agents must receive. Our future work includes the investigation of agent guilt models that capture leakage scenarios.

ACKNOWLEDGEMENT

For all the efforts behind the paper work, we first & foremost would like to express our sincere appreciation to the staff of Dept. of Computer Engg., for their extended help & suggestions at every stage of this paper. It is with a great sense of gratitude that we acknowledge the support, time to time suggestions and highly indebted to our guide **Prof. N. B. Kadu** and **Prof. S. D. Jondhale (H.O.D.)**. Finally, we pay sincere thanks to all those who indirectly and directly helped us towards the successful completion of the paper.

REFERENCES

- [1] P.Buneman, S.Khanna, and W.C.Tan, "Why and Where: A Charaterization of Data provenance," Proc.Eighth Int'l Conf. Database Theory(ICDT '01'),J.V. den Bussche and V.Vianu,eds.,pp.316-330,Jan.2001
- [2] P.Buneman and W.C.Tan,"Provenence in Databases",Proc ACM SIGMOD, pp.1171-1173,2007
- [3] Y.Cui and J.Widom, "Lineage Tracing For General Data Warehouse Transformations," The VLDB J.vol.12,pp.41-58,2003.
- [4] J.J.K.O.Ruanaidh, W.J.Dowling, and F.M.Boland," Watermarking Digital Images For Copyright Protection", IEE Proc.Vision,Signal and Image Processing,vol.143,no.4,pp.250-256,1996.
- [5] F.Hartung and B.Girod,"Watermarking of Uncompressed and Compressed Video," Signal Processing, vol.66, no.3,pp.283-301,1998.
- [6] S.Czerwinski, R.Fromm,and T.Hodes,"Digital Music Distribution and Audio watermarking," <http://www.Scientificcommons.org/43025658>,2007.
- [7] S.Jajodia, P.Samarati, M.L.Sapino,and V.S. Subrahmanian,"Flexible Support For Multiple Access ControlPolicies,"ACMTrans.DatabaseSystems vol.26.no.2.pp.214-260,2001.
- [8] P.Bonatti, S.D.C.di Vimercati,and P.Samarati,"An Algebra For Composing Access Control Policies,"ACM Trans.Information and System Security,vol.5,no.1,pp.1-35,2002.