# Database watermarking, a technological protective measure: Perspective, security analysis and future directions

1 author:

Vidhi Khanduja
Netaji Subhas Institute of Technology
**14** PUBLICATIONS   **138** CITATIONS

SEE PROFILE

# DATABASE WATERMARKING, A TECHNOLOGICAL PROTECTIVE MEASURE: PERSPECTIVE, SECURITY ANALYSIS AND FUTURE DIRECTIONS

Vidhi Khanduja

Department of Computer Engineering, Netaji Subhas Institute of Technology, Delhi, India

vidhikhanduja9@gmail.com

## *Abstract*

Digital Databases dynamically generates a major proportion of the internet content. The databases are created, stored and accessed digitally and transmitted through computer networks. This has grown the potential, sizes and performance of databases in exponential magnitudes. Thus, the need to protect digital databases arises due to the increased vulnerability to copyright and piracy threats originating from the Internet. Both legal and technological measures must be utilized in a synergetic manner to ensure an adequate level of protection. TPMs backed by legal anti-circumvention measures offer a cost-effective solution to database protection. We provide the current state-of-art and analyses of the arena of digital database protection from a combined legal and technical perspective. Our work is more focused on security analysis of the work done so far and providing readers with detailed discussion on the future directions in the domain of digital watermarking of databases.

**Keywords:** Information Security, Digital Watermarking, Right Protection, Tamper Detection.

## 1. INTRODUCTION

Databases are repertoires of knowledge garnered by the collective efforts of mankind through ages and across regions. Digital Databases dynamically generates a major proportion of the internet content. The databases are created, stored and accessed digitally and transmitted through computer networks. This has grown the potential, sizes and performance of databases in exponential magnitudes. Whether they are sold in pieces for data mining applications such as stock market data, consumer behaviour data, power consumption data and weather data or maintained in-house such as product data by e-commerce sites and medical history of patients by hospitals, databases play a pivotal role in all aspects of society.
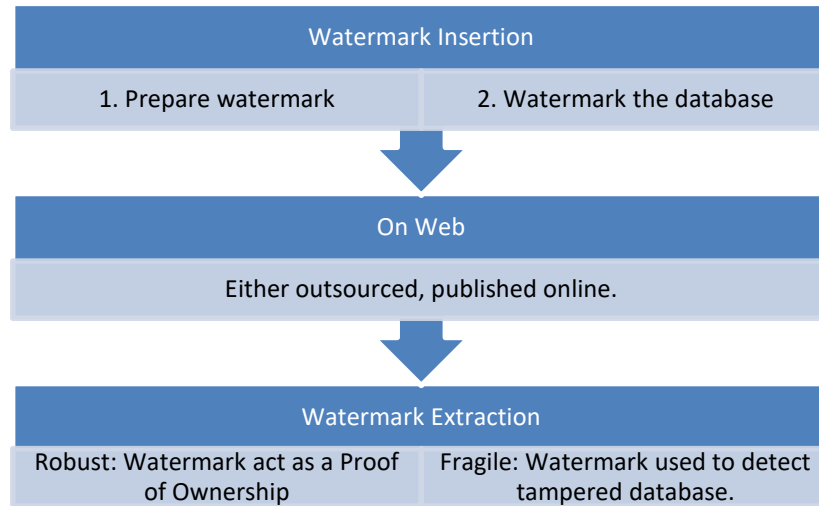
As end users demand more and more information to be available on the net either at low or no cost, database developers are interested in generating revenues from creating databases as this involves intellectual and financial inputs. The need to protect databases arises due to the increased vulnerability to copyright and piracy threats originating from the Internet [1]. Developers have the responsibility to not only supply accurate data, but also ensure its security against illegal copying, hacking or tampering. The digital watermarking of copyrighted works augmented by legal means is fast emerging as an effective and efficient means to protect shared and outsourced databases from infringement.

Section 2 introduces the process of digital watermarking in databases. In section 3, we analyze the arena of digital database protection from a combined legal and technical perspective. Section 4 throws the light on existing work in literature; Section 5 presents the security analyses of the various techniques and Section 6 discusses the future directions. Finally, section 7 concludes the paper.

## 2. AN OVERVIEW OF DIGITAL DATABASE WATERMARKING

Digital watermarking is a viable and cost-effective technological method that protects digital documents such as images, video and databases by *marking* them with some digital pattern. Watermarking algorithms for digital

databases invariably introduce small changes in the data being watermarked with an objective of inserting the mark, but without altering the database in *any significant way*. Watermarking does not completely prevent piracy.



However, it does provide a means to establish the true identity of the owner and deters attempts to plagiarize or distort it [2]. Fig.1 illustrates the process of watermarking databases. The robust watermarking resolves ownership issues while fragile watermarking is for integrity constraints.

**Fig. 1** Watermarking Relational Databases

Unlike encryption and hashing techniques, watermarking does not attempt to hide the data, but instead infuses a kind of ownership proof in the data. Encryption and hashing provide protection to the content by making the information indecipherable by an attacker. On the contrary, digital watermarking works on the principle of modifying the content in a manner so that the usability of the data is retained fully. In fact, an attacker or observer has no way to decipher that a watermark is actually present in the database. However, the watermark does remain within the content inseparably providing a proof of the ownership or a signature to detect tampering or to track the people who may have obtained the content legally and are illegally redistributing it [3].

Many researchers have contributed significantly in the area of watermarking multimedia data such as image, audio and video sources [3-6]. The technique to conceal watermark within databases differs significantly with that of multimedia data [2]. As compared with data in relational databases, multimedia data is voluminous and has a large bandwidth to hide watermarks in a redundant manner. Moreover, the relative spatial positioning of different parts of multimedia data such as image or video is not affected significantly when watermark bits are inserted. Therefore, the quality of image is largely retained.

Relational databases on the other hand, comprise concrete tuples and attributes, each tuple representing a distinct entity. Firstly, we need to disperse the watermark bits across multiple tuples to achieve redundancy. There is no particular ordering between these tuples, so even a subset of tuples can be used. Secondly, embedding of watermark bits invariably introduces perturbations within a database. These perturbations may adversely affect the usability of the database. Therefore, a database watermarking technique must ensure that the usability constraints of the attributes stay within limits when watermark bits are inserted. Usability constraints are the limitations imposed on each attribute. They are decided by the database owner or designer and depend upon its

specific application(s). For example, attribute value should be unique; classification range must remain same before and after concealing watermark etc. [7].

## 3. DATABASE WATERMARKING AS A TECHNOLOGICAL PROTECTION MEASURE WITH LEGAL PROTECTION

The foundation of modern society is intelligible information compiled in innumerable databases. Government departments, corporations, multinational companies, information bureaus and research centers produce databases of government records, medical and legal case records, web pages and collection of literary and artistic works.

Article 1 (2) of Directive 96/9/EC of the European Parliament defines a database as "*a collection of independent works, data or other materials arranged in a systematic and methodical way and individually accessible by electronic or other means*" [8]. According to the Black's Law Dictionary, a database is defined as "*a compilation of information arranged in a systematic way and offering a means of finding specific elements it contains, often today by electronic means*" [9]. The relevance of digital databases is explicitly highlighted in this definition. The distinguishing feature of digital databases is that they can be created, downloaded, value added and digitally re-transmitted with great flexibility and speed.

It requires significant efforts in terms of money, manpower and creative inputs to build high-quality databases. They are thus Intellectual Property in their own right. Given their rich informational content and the ready availability of advanced technologies to communicate and modify them with relative ease, it is imperative to protect digital databases against potential misuse. Technological methods that are designed to protect digital content of any form like text, images and databases are called "Technological Protection Measures (TPM). They include the technologies that can control access to copyrighted digital content or can prevent users from copying such protected content. Watermarking of digital databases is one such TPM that has emerged as an effective means to protect shared and outsourced databases from infringement.

The United States initially followed the "Sweat of the Brow" doctrine to protect databases, which did not require many creative skills such as catalogs and directories. In the landmark case of Feist Publications v. Rural Telephone Service Co., the Supreme Court overturned this and stipulated a modicum of creativity to admit any database under Copyright protection [10].

In Europe before 1996, there was no uniformity in laws regarding protection of databases in member states of the EU. However, with the implementation of European legislation EC Directive 96/9/EC, the treatment of copyrightable databases was harmonized [8]. More significantly, it introduced a new *sui generis* right to bring non-copyrightable databases into the ambit of legal protection.

Focusing on the national scenario, India is a common law country. Databases which have been prepared with cognizable creative efforts in the systematic arrangement of facts are protected under copyright law [11]. However, courts have consistently relied on the "sweat of the brow" doctrine [12]. Some courts have also categorically rejected this doctrine and emphasized on the aspect of creativity in its ad-jurisdiction [13]. India has taken laudable steps to codify significant portions of its traditional knowledge base such as Ayurvedic, Unani, Sidhha etc. by codifying it in the Traditional Knowledge Digital Library (TKDL) [14]. This initiative has succeeded in blocking several false claims of traditional Indian knowledge by individuals and groups to seek patents. World Intellectual

Property Organization (WIPO) and India have ventured together to see how the TKDL model can be emulated by other countries to stop a misappropriation of traditional knowledge [15].

In India, Copyright Amendment Bill 2012, under Section 65A introduces the concept of TPM, as a measure used to enforce restrictions on the use of copyrighted material [16]. Criminal and monetary liabilities are subjected to any person who circumvents the TPMs, with the intention of infringing rights.

Both legal and technological measures must be utilized in a synergetic manner to ensure an adequate level of protection. TPMs backed by legal anti-circumvention measures offer a cost-effective solution to database protection [16-19].

## 4. AN OUTLINE OF DATABASE WATERMARKING RESEARCH

An extensive literature survey is conducted in this domain. In this section we highlight only the major points from contributions exists in literature in the domain of watermarking databases as few survey papers already exists in the literature [20-23]. However, they only discuss their techniques and classify them accordingly. We provide the lacunae and advantages of primary research in this domain. Additionally, our work is more focused on security analysis of the work done so far and providing readers with detailed discussion on the future directions. Fig. 2 and Fig. 3 illustrate the features of watermarking techniques. The watermark is either randomly selected binary stream, or Object Identifier selected by owner or based on owner's biometric trait. The watermark can be embedded in all tuples (AT) or few selected tuples (MT). Similarly, within a tuple, single (SA) or multiple attributes (MA) can be selected for embedding watermark bit. Finally, within an attribute either single bit (SA) is modified or multiple bits (MB) depending on various watermarking approach [21].
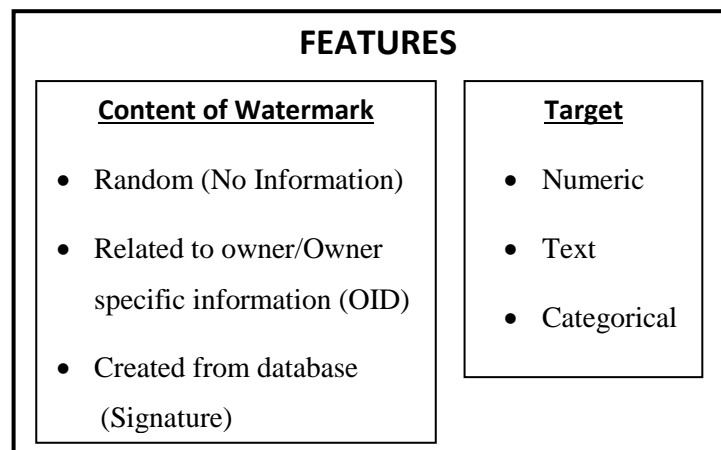


**FEATURES**

| **Content of Watermark** | **Target** |
|---|---|
| • Random (No Information) | • Numeric |
| • Related to owner/Owner specific information (OID) | • Text |
| • Created from database (Signature) | • Categorical |

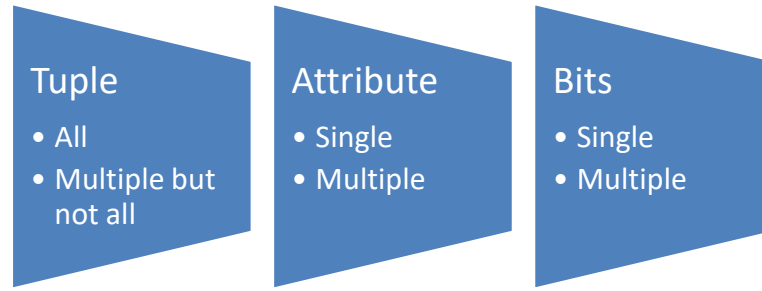**Fig. 2** Features of watermarking techniques

**Fig. 3** Various Granularity levels

Table 1 and 2 shows the comparison of various existing watermarking techniques targeting numeric as well as other attributes. Table 3 enlists various fragile watermarking techniques. These techniques are able to detect any temperedness introduced in the dataset. The watermark acts as a signature and is prepared from the dataset. It is either embedded making a distortion in the dataset, or saved with trusted third party making a distortion-free technique. The localization of position where error has occurred varies with technique.

Based on our analysis and past researchers work; we found that most of the work is focused on traditional relational databases. The existing techniques have served the basic purpose of technological protection of digital content well. But they need to be further improved by increasing the degree of robustness against malicious attacks, minimizing the distortions produced due to the process of watermarking itself and enhancing the reliability with which the proof of ownership is established in case of conflicts.

**Table 1: Summary and comparison of robust watermarking techniques targeting numeric attributes.**

| Proposed Schemes | Watermark Content | Features | Weakness |
|---|---|---|---|
| Agrawal *et. al.* [2][24] | Random | 1. Identified need of watermarking Relational databases<br>2. Single attribute and single bit (SASB) selection approach. | 1. Single bit-watermark<br>2. Primary Key (P.K) dependent technique<br>3. Not resilient to Subset deletion, attribute re-order Attack |
| Xinchun *et.al.* [25] | OID | 1. Weighted algorithm based attribute selection<br>2. SASB approach. | 1. Prone to Bit- attacks<br>2. P.K dependent technique |
| Farfoura *et.al.* [26] | OID | 1. Reversible technique<br>2. Use of Time stamping protocol<br>3. Single attribute and multiple bit (SAMB) selection | 1. P.K dependent technique<br>2. Not resilient to attribute Re- order attack<br>3. Linear Transformation attack not addressed |
| Zhou *et. al.*[27] | Image | 1. Random changes not visible<br>2. Use of TTP<br>3. SASB selection approach. | 1. Not resilient to Subset deletion attack<br>2. Prone to Bit-attacks |

| | | | |
|---|---|---|---|
| Sun *et.al.*[28] | Image | 1. Multiple images are embedded<br><br>2. SASB selection approach. | 1. P.K dependent technique<br>2. Recognizable Common pattern /correlation in selection (Same value used for tuple, attribute and bit selection.) |
| Wang *et. al.* [29] | Biometric trait | 1. Use of Owner's voice as watermark<br><br>2. SASB selection approach. | 3. Biometrics is not systematically addressed<br>2. Recognizable Common pattern in selection (Same parameter used for tuple, attribute & bit selection.)<br>3. No technique outlined for handling false claims of ownership |
| Sion *et. al.* [7] | Random | 1. Partition Statistics are used<br>2. Usability Constraints considered<br>3. P.K independent<br>4. SAMB selection | 1. Prone to Synchronization errors as marker tuples are used<br>2. Not resilient to Subset Deletion and Alteration attack<br>3. No clear systematic approach for data manipulation |
| Shehab *et.al.* [30] | Random | 1. Partitioning without marker tuples<br>2. Optimization based technique.<br>3. Multiple attributes and multiple bits (MAMB) selection | 1. Computationally not efficient<br>2. Linear Transformation attack not addressed<br>3. P.K dependent technique |
| Khanduja *et.al.* [31] | OID | 1. Optimization based technique.<br>2. Use of Bacterial Foraging algorithm<br>3. MAMB selection | 1. Technique applicable to only Numeric attributes. |
| Iftikar *et.al.*[32] | Optimal watermark prepared from database information | 1. Reversible watermarking technique<br><br>2. Optimal watermark creation through Genetic Algorithm<br><br>3. All tuples, MAMB selection | 1. Semi-blind technique.<br><br>2. Not resilient to attribute deletion attack. |

**Table 2: Summary and comparison of robust watermarking techniques targeting other data types.**

| Proposed Schemes | Target | Water-mark | Features | Weakness |
|---|---|---|---|---|
| Odeh *et. al.*[33] | Time | Image | 1. Large bit capacity for watermark embedding as SAMB selection<br>2. No effect on usability of data | 1. Not resilient to Subset deletion attack as Marker tuples are used<br>2. Not resilient to Subset Alteration attack |
| Al-Haj *et.al.*[34] | Text (Non-Numeric) | Image | 1. Large Bit capacity for watermark insertion | 1. Not resilient to Subset alteration attack as double spaces are added |

| | | | 2. Extra spaces embedded<br>3. SAMB selection | 2. Security issues not addresses like secret key not used |
|---|---|---|---|---|
| Hanyurwimfura *et.al.* [35] | Non-numeric | Random bit-stream | 1. Watermark bit is embedded in multiword attribute with the lowest Levenshtein Distance<br>2. Large bit capacity as MAMB selection | 1. No discussion on watermark<br>2. P.K dependent technique<br>3. Not so secure technique as is easily detectable |
| Sion *et. al.*[36-37] | Robust | Bit string | 1. Identified the need of watermarking categorical data<br>2. Error Correcting Code is applied<br>3. SAMB selection | 1. Introduce distortions in categorical data<br>2. P.K dependent technique<br>3. Less Resilient against malicious attacks |

**Table 3: Summary and comparison of fragile watermarking techniques.**

| Proposed Schemes | Watermark | Data Type | Features | Weakness |
|---|---|---|---|---|
| Li *et. al.*[38] | Signature | Categorical | 1. Distortion free technique<br>2. Localization up to group level | 1. Not resilient to Tuple re-ordering attack<br>2. P.K dependent technique |
| Guo *et.al.* [39] | Signature | Numeric | 1. Partitioning based technique<br>2. Localization up to tuple level | 1. Perturbations in LSBs not detectable<br>2. P.K dependent technique |
| Khan *et.al.* [40] | Signature | Categorical | Distortion free technique | 1. P.K dependent technique<br>2. Prone to Attribute-value substitution attack |
| Khataeimara gheh *et.al.* [41] | Signature | Numeric | 1. Localization of the tampered data<br>2. Technique to recover data | 1. Fails to localize for two or more than two updations in different tuples/ attributes<br>2. Fails to recover when an attribute is deleted.<br>3. No focus on information recovery |
| L.Camara, *et.al.* [42] | Signature | Watermark not embedded in database | Distortion-free technique | 1. No discussion on watermark preparation<br>2. P.K dependent technique |
| Kamel [43] | Secret number | Indirect embedding | Distortion-free technique based on reordering of tuples | 1. Not resilient to tuple alteration attack<br>2. No discussion on watermark preparation. |

## 5. SECURITY ANALYSIS

The techniques proposed numerous processes of concealing watermarks. In most of them, the watermarks are concealed within a database. We now analyse the security of such techniques considering generalized

situation. The analysis on various attacks is presented in [21]. In this work, we analyse an important aspect of security by calculating the probability to find the potential locations. If an attacker, say Mallory, tries to destroy the watermark to claim the database to be hers. Under such circumstances, we analyse the difficulty level of Mallory to find the positions of embedded bits to alter the watermark. Another important concern is False Hit Rate. Any watermarking model is considered to more robust if its False Hit Rate is minimum. Section 5.1 analyses the security by calculating the probability to find potential locations and in Section 5.2 below we discuss the False Hit Rate for various watermarking techniques.

### 5.1. Probability to find Potential locations

We calculate the probability to find the potential locations based on granularity levels of various embedding procedures as illustrated in fig. 3 and table 1 and 2. The watermark is converted to binary and then bit/s of watermark is embedded in the attribute value (converted to binary form). Let the number of bits in a watermark be $L_w$ . The process requires maximum of $L_w$ attributes to completely embed watermark once considering a single bit, concealed within an attribute. Positions where watermark is to be concealed is securely selected (i.e. using secret parameters). Let $N_a$ be number of attributes in a database and $N_{pa}$ is number of permissible attributes where watermark can be concealed such that $N_{pa} \leq N_a$. $N_{pa}$ is decided by the owner of the database considering usability constraints. Let, the attribute length be $L_a$ and number of permissible bits within $L_a$, where watermark bit can be embedded without violating the usability constraint of an attribute is $L_{pa}$ such that $L_{pa} \leq L_a$. Table 4 enlists the other symbols used herewith.

Table 4: Notations

| Symbol | Meaning /Explanation |
|---|---|
| $L_w$ | Length of watermark (in bits) |
| $L_a$ | Length of an attribute (in bits) |
| $L_{pa}$ | Number of permissible bits out of $L_a$, in an attribute where watermark can be embedded. |
| $N_a$ | Number of attributes in a database |
| $N_{pa}$ | Number of permissible attributes out of $N_a$ where watermark can be embedded. |
| $N_t$ | Number of tuples in a database |
| $M_a$ | Multiple attributes selected out of $M_{pa}$; attributes where embedding occurs. $M_a \leq N_{pa}$ |
| $M_t$ | Multiple tuples out of $N_t$ where watermark is embedded, $M_t \leq N_t$. |
| $M_b$ | Multiple bits within a selected attribute to conceal watermark bits, $M_b \leq L_{pa}$. |

Let us first calculate the probability to correctly select target location within a single tuple.

The probability to correctly select single watermarked attribute out of $N_{pa}$ is given by

$$P_{SA} = 1/N_{pa} \tag{1}$$

The probability to correctly select bit position within a selected attribute $P_{SB} = 1/L_{pa}$ \hfill (2)

In certain techniques, multiple attributes are selected per tuple [44-45]. In such cases, probability to select multiple attributes say, $M_a$ attributes out of $N_{pa}$

$$P_{MA} = 1/\left(\begin{smallmatrix}N_{pa}\\M_a\end{smallmatrix}C\right) \tag{3}$$

In certain techniques, multiples bits positions are selected within a selected attribute to embed multiple bits say $M_b$, of watermark [44]. In such cases, the probability to correctly select $M_b$ embedding locations is given by:

$$P_{MB} = 1/\left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right) \tag{4}$$

Total probability to correctly choose single watermarked location within a single earmarked attribute of a tuple is calculated using Eq(1) and Eq(2) as

$$P_{SASB} = 1/\left(N_{pa} * L_{pa}\right) \tag{5}$$

Now, considering entire database comprising $N_t$ tuples, the probability in Eq (5) becomes:

$$P_{ATSASB} = \left(1/(N_{pa} * L_{pa})\right)^{N_t} \tag{6}$$

This is the probability when all the tuples are used to embed the watermark [30, 34]. This produces more alterations in the database at the cost of more potential locations for concealing watermark bits [21].

Let us suppose watermark is not embedded in all the tuples. Certain tuples say, $M_t$ is selected for embedding watermarks using owner selected secret parameters [2, 24-26, 29, 45]. We now calculate probability in each of the following cases:

i.  **Case MTSASB:** The probability to correctly select $M_t$ tuples out of $N_t$ (considering single bit and single attribute selection Procedure) using Eq (1) and (2) is given by:

$$P_{MTSASB} = \left\{\frac{M_t * N_t}{N_t * N_{pa} * L_{pa}} * \frac{(M_t * N_t) - 1}{(N_t - 1)N_{pa} * L_{pa}} * \frac{(M_t * N_t) - 2}{(N_t - 2)N_{pa} * L_{pa}} * \dots \right.$$
$$\left. * \frac{1}{(N_t - M_t * N_t + 1)N_{pa} * L_{pa}}\right\}$$
$$= \left(1/\left(\begin{smallmatrix}N_t\\M_t\end{smallmatrix}C\right)\right) * \left(1/\left(N_{pa} * L_{pa}\right)\right)^{M_t} \tag{7}$$

ii.  **Case MTSAMB:** The probability to correctly select $M_t$ tuples out of $N_t$ (considering multiple bits and single attribute selection Procedure) using Eq (1) and (4) is given by: P$_{MTSAMB}$

$$P_{MTSAMB} = \left\{\frac{M_t * N_t}{N_t * N_{pa} * \left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right)} * \frac{(M_t * N_t) - 1}{(N_t - 1)N_{pa} * \left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right)} * \frac{(M_t * N_t) - 2}{(N_t - 2)N_{pa} * \left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right)} * \dots \right.$$
$$\left. * \frac{1}{(N_t - Mt * N_t + 1)N_{pa} * \left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right)}\right\}$$
$$= \left(1/\left(\begin{smallmatrix}N_t\\M_t\end{smallmatrix}C\right)\right) * \left(1/\left(N_{pa} * \left(\begin{smallmatrix}L_{pa}\\M_b\end{smallmatrix}C\right)\right)\right)^{M_t} \tag{8}$$

iii.  **Case MTMAMB:** Total probability to correctly select all the watermarked locations considering multiple bits per attribute and multiple attribute selection within $M_t$ selected tuples is given by:

$$P_{MTMAMB} = \left\{ \frac{M_t * N_t}{N_t * \binom{N_{pa}}{M_a}C * \binom{L_{pa}}{M_b}C} * \frac{(M_t * N_t) - 1}{(N_t - 1) * \binom{N_{pa}}{M_a}C * \binom{L_{pa}}{M_b}C} * \frac{(M_t * N_t) - 2}{(N_t - 2) * \binom{N_{pa}}{M_a}C * \binom{L_{pa}}{M_b}C} \right.$$

$$\left. * \dots \dots * \frac{1}{(N_t - M_t * N_t + 1) * \binom{N_{pa}}{M_a}C * \binom{L_{pa}}{M_b}C} \right\}$$

$$= \left\{ \frac{1}{\binom{N_t}{M_t}C} \right\} * \left\{ \frac{1}{\binom{N_{pa}}{M_a}C * \binom{L_{pa}}{M_b}C} \right\}^{Mt} \tag{9}$$

Probability to detect potential locations in all the four above cases i.e. ATSASB, MTSASB, MTSAMB and MTMAMB is calculated and values are plotted fig. 4. The graph shows that the probability to find potential locations by an attacker without the knowledge of secret parameters decreases as we move from Eq. (7) for MTSASB to Eq.(8) for MTSAMB and then (9) for MTMAMB and finally we get least probability in case of Eq. (6) for ATSASB. The simulation parameters taken are enlisted in table 5.

Table 5: Simulation values

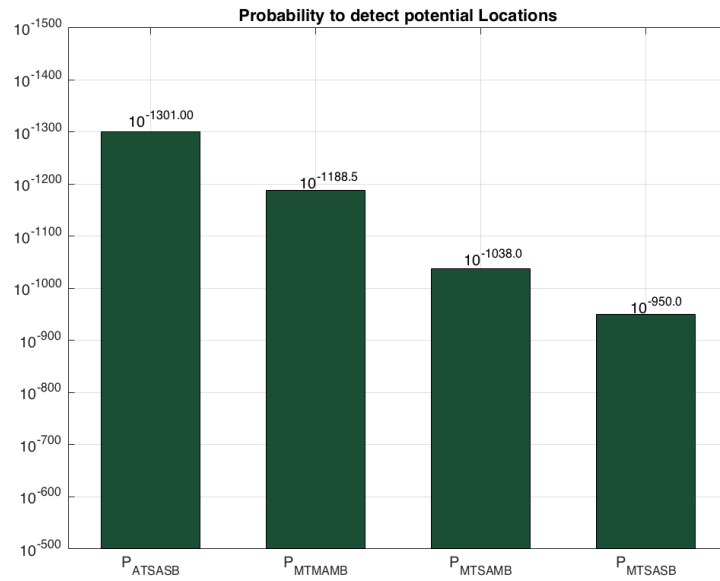| Simulation Parameters | Values |
|---|---|
| $N_t$ | 1000 |
| $N_{pa}$ | 5 (Varied from 3 to 7 ) |
| $M_a$ | 2 |
| $L_{pa}$ | 4 (Varied from 2 to 8) |
| $M_b$ | 2 |
| $M_t$ | 50% of $N_t$ (Varied from 30% to 80%) |



**Fig. 4** Probability to select potential locations in case of ATSASB, MTSASB, MTSAMB and MTMAMB.

Probability to detect potential locations for single attribute, single bit selection in all tuples of the database is calculated using Eq(6). The variations in probability is recorded in table 6 by varying $N_t$, $N_{pa}$ and $L_{pa}$. As number of tuples increases, number of potential location increases thereby, decreasing the probability. Table 7 contains the probability in case of MTSASB. The values are recorded at different values of $M_t$.

Considering $M_t$=50% of $N_t$ we also plotted the graph showing variation in probability calculated using Eq. 6, 7, 8 and 9 respectively shown in fig. 5. It shows the decrease in probability to detect potential locations with the increase in $N_t$. Increase in the number of tuples causes increase in the number of potential locations thus; probability to detect them correctly further decreases in all four cases.

Table 6 : Probability to detect potential locations for single attribute, single bit selection in all tuples of the database by varying $N_t$, $N_{pa}$ and $L_{pa}$

| P ATSASB | $N_t$=100 | $N_t$=500 | $N_t$=1000 | $N_t$=1500 | $N_t$=5000 |
|---|---|---|---|---|---|
| P ATSASB($N_{pa} = 7 \ L_{pa} =4$) | 1.9E-145 | 2.6E-724 | 6.9E-1448 | 1.8E-2171 | 1.6E-7236 |
| P$_{ATSASB}$ ($N_{pa} = 5 \ L_{pa}$ =4) | 7.8E-131 | 3.0E-651 | 9.3E-1302 | 2.8E-1952 | 7.0E-6506 |
| P$_{ATSASB}$($N_{pa} = 3 \ L_{pa}$ =4) | 1.2E-108 | 2.5E-540 | 6.5E-1080 | 1.6E-1619 | 1.2E-5396 |
| P$_{ATSASB}$ ($N_{pa} = 5 \ L_{pa}$ =8 | 6.2E-161 | 9.3E-802 | 8.7E-1603 | 8.1E-2404 | 5.0E-8011 |
| P$_{ATSASB}$ ($N_{pa} = 5 \ L_{pa}$ =2) | 1.0E-100 | 1.0E-500 | 1.0E-1000 | 1.1E-1500 | 1.0E-5000 |

Table 7: Probability to detect potential locations for single attribute, single bit selection in selected tuples of the database by varying $N_t$, $N_{pa}$, $L_{pa}$ and a) $M_t$=80% b) $M_t$=50% c) $M_t$=30%

| MTSASB, $M_t$=80% | $N_t$=100 $M_t = 80$ | $N_t = 500$ $M_t = 400$ | $N_t = 1000$ $M_t = 800$ | $N_t = 1500$ $M_t = 1200$ | $N_t$=5000 $M_t = 4000$ |
|---|---|---|---|---|---|
| P$_{MTSASB}$($N_{pa} = 7, L_{pa} = 8$) | 2.6E-161 | 2.2E-807 | 4.2E-1615 | 6.2E-2423 | 3.1E-8078 |
| P$_{MTSASB}$ ($N_{pa} = 5 \ L_{pa}$ =4) | 1.5E-125 | 1.6E-628 | 2.2E-1257 | 2.4E-1886 | 1.3E-6289 |
| P$_{MTSASB}$ ($N_{pa} = 5, L_{pa} = 8$ | 1.2E-149 | 6.2E-749 | 3.4E-1498 | 1.4E-2247 | 1.0E-7493 |
| P$_{MTSASB}$ ($N_{pa} = 5, L_{pa} = 2$) | 1.8E-101 | 4.1E-508 | 1.5E-1016 | 4.1E-1525 | 1.7E-5085 |
| P$_{MTSASB}$ ($N_{pa} = 3, L_{pa} = 4$) | 8.7E-108 | 8.8E-540 | 6.8E-1080 | 4.0E-1620 | 3.3E-5402 |
| P$_{MTSASB}$ ($N_{pa} = 7, L_{pa} = 4$) | 3.1E-137 | 5.7E-687 | 2.8E-1374 | 1.0E-2061 | 4.0E-6874 |

b) Considering $M_t$=50%

| MTSASB, $M_t$=50% | $N_t$=100 $M_t = 50$ | $N_t = 500$ $M_t = 250$ | $N_t = 1000$ $M_t = 500$ | $N_t = 1500$ $M_t = 750$ | $N_t$=5000 $M_t = 2500$ |
|---|---|---|---|---|---|
| P$_{MTSASB}$($N_{pa} = 7, L_{pa} = 8$) | 3.8E-117 | 5.6E-587 | 2.9E-1174 | 1.0E1761 | 2.2E-5874 |
| P$_{MTSASB}$ ($N_{pa} = 5 \ L_{pa}$ =4) | 8.8E-95 | 3.4E-475 | 1.1E-950 | 2.3E-1426 | 1.7E-4756 |
| P$_{MTSASB}$ ($N_{pa} = 5, L_{pa} = 8$ | 7.8E-110 | 1.9E-550 | 3.4E-1101 | 3.9E-1652 | 4.7E-5509 |
| P$_{MTSASB}$ ($N_{pa} = 5, L_{pa} = 2$) | 1.0E-79 | 6.2E-400 | 3.7E-800 | 1.3E-1200 | 6.6E-4004 |
| P$_{MTSASB}$ ($N_{pa} = 3, L_{pa} = 4$) | 1.0E-83 | 1.0E-419 | 9.5E-840 | 5.7E-1260 | 7.4E-4202 |
| P$_{MTSASB}$ ($N_{pa} = 7, L_{pa} = 4$) | 4.3E-102 | 1.0E-511 | 9.7E-1024 | 5.9E-1536 | 8.4E-5122 |

c) Considering $M_t$=30%

| MTSASB, $M_t$=30% | $N_t$=100 $M_t$ =30 | $N_t = 500$ $M_t = 150$ | $N_t = 1000$ $M_t = 300$ | $N_t = 1500$ $M_t = 450$ | $N_t$=5000 $M_t = 1500$ |
|---|---|---|---|---|---|

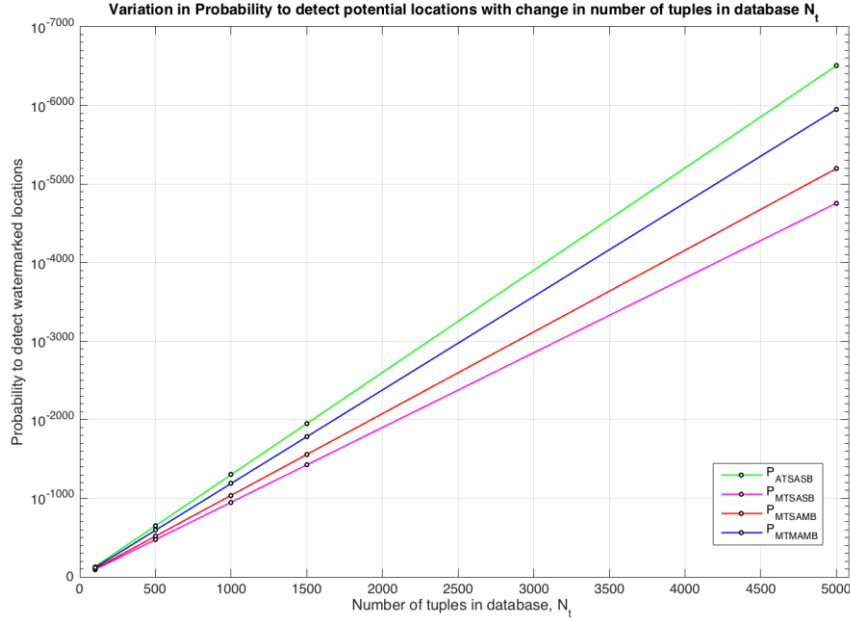| | | | | | |
|---|---|---|---|---|---|
| $P_{MTSASB}(N_{pa} = 7, L_{pa} = 8)$ | 1.2E-78 | 3.4E-394 | 6.4E-789 | 1.0E-1183 | 1.4E-3947 |
| $P_{MTSASB}(N_{pa} = 5\ L_{pa} = 4)$ | 3.2E-65 | 4.1E-327 | 9.0E-655 | 1.8E-982 | 7.9E-3277 |
| $P_{MTSASB}(N_{pa} = 5, L_{pa} = 8)$ | 2.9E-74 | 2.8E-372 | 4.4E-745 | 6.2E-1118 | 2.2E-3728 |
| $P_{MTSASB}(N_{pa} = 5, L_{pa} = 2)$ | 3.4E-56 | 5.8E-282 | 1.8E-564 | 5.2E-847 | 2.7E-2825 |
| $P_{MTSASB}(N_{pa} = 3, L_{pa} = 4)$ | 1.4E-58 | 7.8E-294 | 3.2E-588 | 1.2E-882 | 4.6E-2944 |
| $P_{MTSASB}(N_{pa} = 7, L_{pa} = 4)$ | 1.3E-69 | 4.9E-349 | 1.3E-698 | 3.1E-1048 | 5.0E-3496 |



**Fig. 5** Effect of total number of tuples on probability to detect watermarked locations

Similarly, the effect of $M_t$ on probability is plotted in fig. 6. As number of marked tuples, $M_t$ increases, the number of potential location increases, thereby decreasing the probability in all the three cases. However, MTMAMB is more dependent on $M_t$ as compared to MTSASB. Probability in case of ATSASB does not depend on $M_t$. It was also observed that this probability also depend on $N_{pa}$. Fig. 7 depicts the effect of $N_{pa}$ on the probability. The value of $N_{pa}$ is varied from 3 to 7 and probability is observed in each of the four cases. Probability in case of MTMAMB shows maximum dependency on $N_{pa}$, as it shows variation/decrease in probability.

Next, the effect of $L_{pa}$ on probability is observed and plotted in fig. 8. The value of $L_{pa}$ is varied from 2 to 8 and the respective probability is recorded in each of the four cases. Variation in probability is almost equal in case of MTSAMB and MTMAMB showing equal dependency which can be verified from Eq(8) and Eq(9).
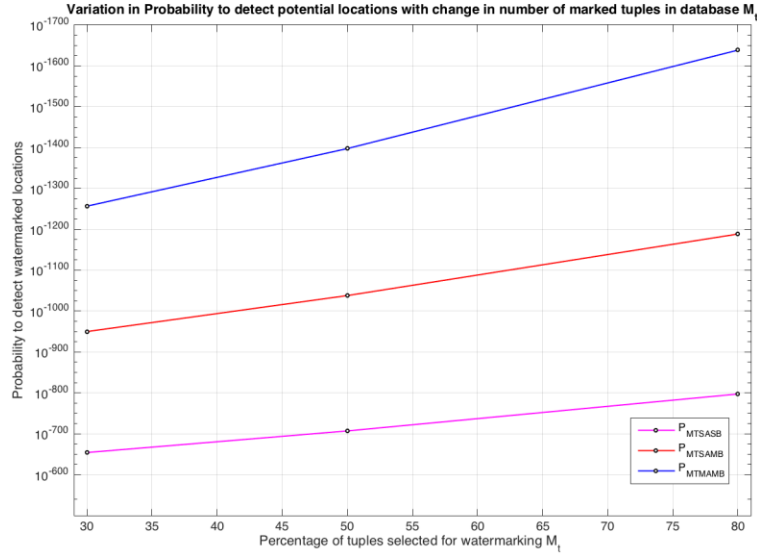
**Fig. 6** Effect of total number of marked tuples on probability to detect watermarked locations
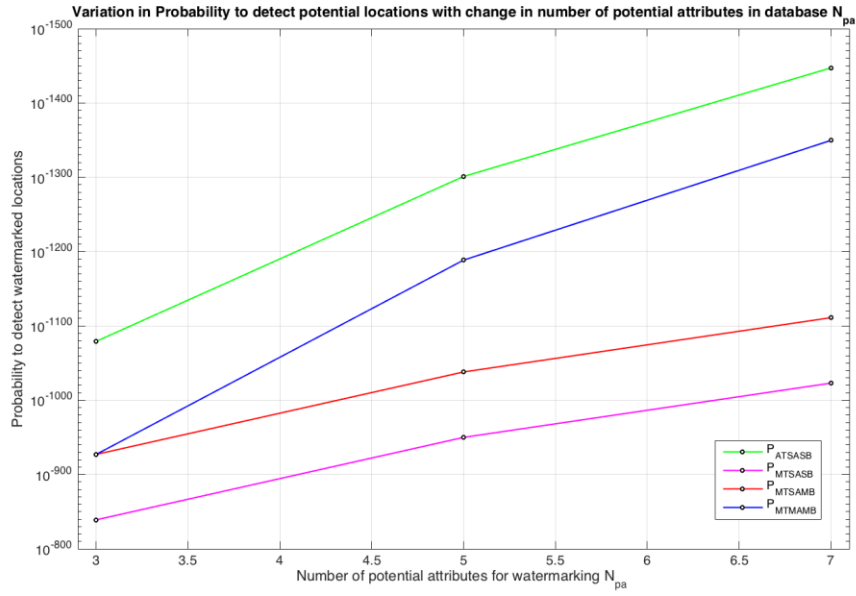


**Fig. 7** Effect of total number of potential attributes on probability to detect watermarked locations
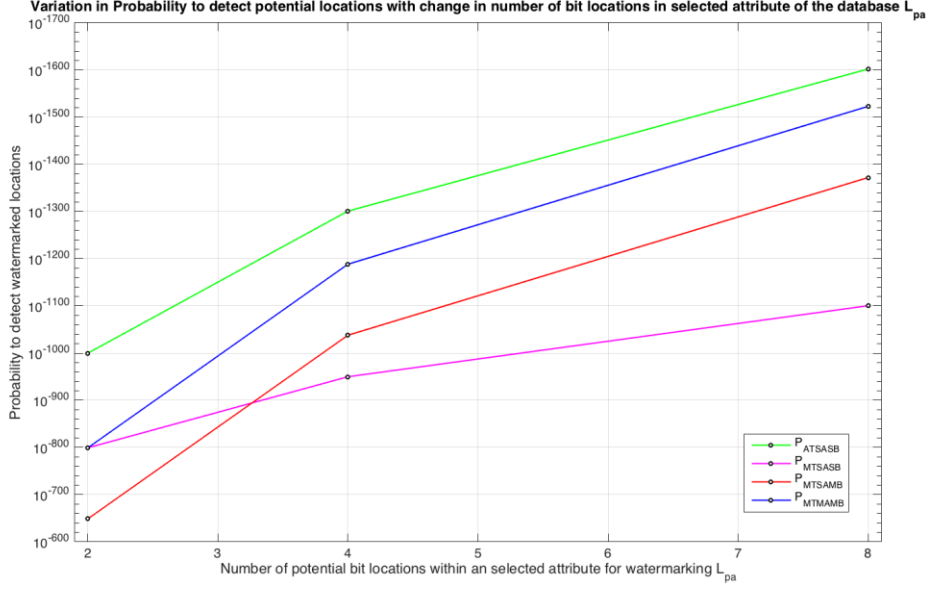
Fig. 8 Effect of total number of permissible bits within a selected attribute on probability to detect watermarked locations

### 5.2. False Hit Rate

We now calculate the False Hit Rate $P_{FH}$ to analyse the security. $P_{FH}$ is the probability of extracting a valid watermark from a non-watermarked database. Our $P_{FH}$ analysis is grounded upon the approach followed in [2]. In [2] each watermark bit is embedded once. However, all the recent work embeds single watermark bit repeatedly to make the system more robust against malicious attacks such as subset alteration, insertion and deletion attacks. The final watermark bit is decided using majority voting among multiple extracted bits. In this work, we extend the calculation done in [2] for repetitive embedding of a watermark bit.

Let a watermark bit $b_i$ be embedded $N_r$ times in a database. Each bit $b_i^*$ is extracted from a database with the probability of $1/2$ to match with the original embedded watermark bit. The final watermark bit is decided based on majority voting. We define $P_{bit}$ as the probability of correctly extracting one watermark bit after majority voting by sheer chance. This is equivalent to saying that $P_{bit}$ is the probability that at-least one more than half of $N_r$ bits can be detected from non-watermarked relation by sheer chance. We now calculate $P_{bit}$ using binomial distribution considering $N_r$ independent trials [46].

$$P_{bit} = \sum_{j=N_r/2+1}^{N_r} b\left(j; N_r, \frac{1}{2}\right) = B\big((0.5 * N_r + 1); N_r, 0.5\big) \tag{10}$$

Where, $b(j; n, p)$ be the probability of obtaining $j$ success in $n$ Bernoulli trials with probability $p$ for success and $1 - p$ for failure given by equation (11).

$$b(k; n, p) = \frac{n}{k}C * p^k * (1 - p)^{n-k} \tag{11}$$

We use $B(k; n, p)$ referred as cumulative binomial probability representing the probability of obtaining at-least $k$ success from $n$ Bernoulli's trials.

$$B(k; n, p) = \sum_{i=k}^{n} b(i; n, p) \tag{12}$$

For a watermark of length $N_w$, the false hit rate $P_{FH}$ is given by:

$$P_{FH} = B(\tau_w * N_w; N_w, P_{bit}) \tag{13}$$

Where, $N_w$ is length of watermark in bits and $\tau_w$ is watermark threshold such that if $\tau_w * N_w$ watermark bits correctly matches the original watermark out of $N_w$ bits; we can claim that the suspected database is ours. The value of $N_r$ will vary for different embedding schemes as follows.

$$\text{In case of ATSASB, } N_r = \left\lfloor \frac{N_t}{N_w} \right\rfloor. \tag{14}$$

$$\text{For MTSASB, } N_r = \left\lfloor \frac{M_t}{N_w} \right\rfloor. \tag{15}$$

$$\text{For, MTSAMB, } N_r = \left\lfloor \frac{M_t * M_b}{N_w} \right\rfloor. \tag{16}$$

$$\text{And in case of MTMAMB, } N_r = \left\lfloor \frac{M_t * M_a * M_b}{N_w} \right\rfloor \tag{17}$$

We took $N_t = 1000, M_b = 2, M_t = 500, M_a = 2, N_w = 50$ and plotted the graph by varying $\tau_w$. Fig. 9 illustrates that false hit rate is monotonically decreasing with increase in threshold $\tau_w$. We observed the effect of $N_r$ on $P_{FH}$ and plotted the values in same graph. It is revealed that $P_{FH}$ increases with increase in $N_r$. Table 7 shows the values obtained. $N_r$ is maximum in case of MTMAMB, resulting in maximum redundancy and thus $P_{FH}$ also gets maximum value. The watermarking model is considered best , if it has minimum False Hit Rate. False Hit Rate can be decreased by lowering $N_r$.

Table 7. Values of $N_r$ and False Hit Rate

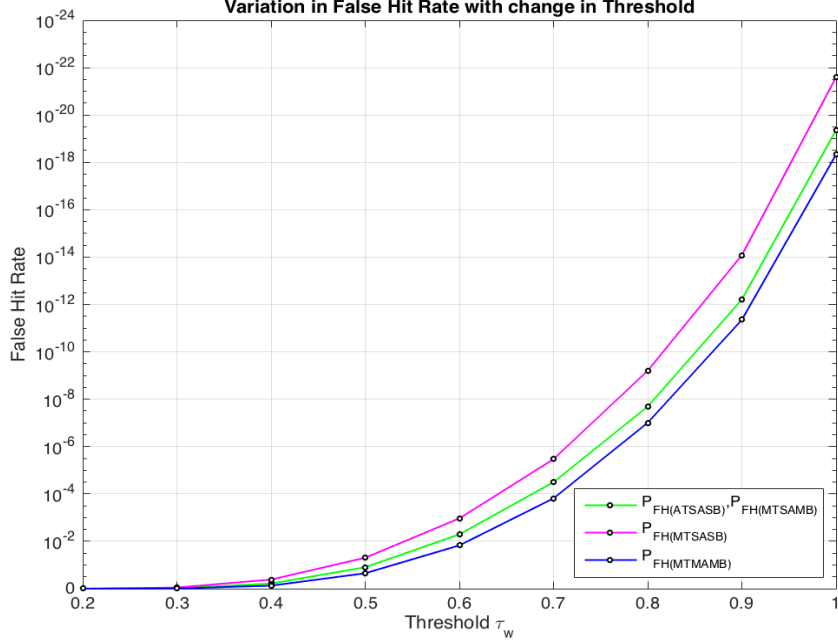|  | $N_r$ | $P_{FH}$ at $\tau_w$=0.6 |
|---|---|---|
| ATSASB | 20 | 0.0051 |
| MTSASB | 10 | 0.00108 |
| MTSAMB | 20 | 0.0051 |
| MTMAMB | 40 | 0.0149 |

**Fig. 9** Effect of threshold $\tau_w$ on False Hit Rate $P_{FH}$ for all four watermarking approach.

## 6. FUTURE DIRECTIONS

Databases that populate the web need reliable technological measures for their protection. We have identified four security aspects that need powerful combative measures for protection.

**(i)** **Ownership Proof**: To provide evidence in ownership disputes, the watermark is securely concealed within the database. Concealed watermark should be robust to various malicious attacks.

**(ii)** **Tamper Detection:** To detect the occurrence of perturbations in the database, watermark is concealed within it. Embedded watermark should be fragile to various malicious attacks.

**(iii)** **Information Recovery:** The information contained within a database must be protected. The watermark must contain the recovery information to restore the lost information from a database.

**(iv)** **Authentication:** To prove the authenticity of the sender, an identity of a sender/owner is embedded as a watermark within the database. This is achieved by embedding biometric trait of the owner as a watermark.

Significant work is done in the field of digital watermarking of databases for protection of first two security aspects; ownership proof and tamper detection [24-40]. We observed that the statistical watermarking techniques need to be improved for better applicability [7, 30, 31]. We have identified lacunae in existing technique and enumerated them in table 1, 2, and 3. We found that it is worth examining these techniques for the challenging application of ownership proofs. This has given us an objective of increasing the reliability of existing robust watermarking technique.

There is an ample opportunity to explore new digital watermarking techniques that cater to the requirements for emerging applications. We found that a watermark can be fruitfully applied as an information carrier to address

the other two security issues namely, information recovery and reliable authentication of content providers. Therefore, we venture to address these additional security aspects.

Some databases such as meteorological department, medicine, military, transportation etc. contain critical attributes. These high risk databases need to be preserved. Moreover, the information in the critical attribute is subject to various malicious attacks resulting in threat to their ownership. This motivated to highlight the need to develop a scheme that recovers lost information and also resolves ownership issues.

A fragile watermarking technique for recovering data is proposed by Khataeimaragheh *et.al.* [41]. The proposed scheme can detect and correct perturbations in RDBs. However, it was found that the probability of accurately rectifying errors reduces drastically as the number of errors exceeds two in [41]. In another work in literature, proposed technique deciphers data in terms of information it represents [44]. It recovers the information lost from altered as well as deleted data irrespective of number of errors occurred in database [44].

We move ahead in same direction, highlighting the utility of watermark as carrier of the information to address another security issue, i.e. to authenticate the owner or content provider of a database. Scant attention is given towards exploring the potential of biometrics to attain the ownership in case of digital databases. Wang *et.al.* suggested a use of owner's voice to establish ownership [29]. This technique was improvised in [45]. In [45] a robust technique that protects the ownership rights with a high degree of accuracy is proposed. Their technique proposes embedding watermarks in multiple attributes thus, enabling better resilience in comparison with [29]. The relative evaluation approach is applied by comparing scores of voice samples of the contenders; thus identification of the owner is achieved even if a watermark extracted is of degraded quality. Thrust lies in investigating multiple biometrics for improvising the reliable ownership proof in noisy environments.

Further, as we move towards distributed, real time and mobile applications, the time taken to embed a watermark must be reduced. The challenge is to consider the protection of emerging web databases such as object oriented and object-relational web databases that have started gaining popularity and are outsourced or shared on the web. Further dynamic web databases such as NoSql databases are being widely adopted for large-scale data applications and social networking sites. Since the traditional database systems weren't designed or optimized for such enormous size of data with complex data types and structures, the need of Not Only SQL (NoSQL) databases increased. These databases may or may not have features of traditional relational databases in favour of better horizontal scaling facilities, which is a problem for RDMS or having schema-less document based objects, which allow capturing complex structures such as data coming from several different sensors and also allowing faster access in some cases. They need urgent protection against piracy and integrity losses as it seemed to have escaped the attention of researchers.

## 7. CONCLUSION

The area of digital watermarking is rife with challenges, and ample research is still ongoing. We emphasized granularity level and thus, categorized watermarking techniques into four types, i) ATSASB ii) MTSASB iii) MTSAMB and iv) MTMAMB. We have analysed the security of watermarking techniques w.r.t. ability to locate watermarked positions within a database without knowledge of secret parameters. We have theoretically analysed its dependency on various parameters; i) $N_r$ ii) $M_t$ iii) $L_{pa}$ iv) $N_{pa}$. Increase in the number of tuples causes increase in the number of potential locations thus; probability to detect them correctly further decreases in all four cases. For a secure system, this probability should be least. We further

calculated the false hit rate of the watermarking techniques based on their categories. More the redundancy, more the false hit rate and less robust the approach.

Our research shows the multiple directions which are worth addressing in future endeavours.

**1.** There is a need to design suitable watermarking techniques targeting emerging web databases to provide protection against various security aspects.

**2.** Future work can focus on exploiting watermarking for extracting and thus, regenerating the lost/perturbed information concealed within a database.

3. Extending the use of watermarking to address additional security issues such as authentication and information recovery.

Given these promising research directions, the domain of watermarking can be further enriched with a thrust on database security.

### *References*

[1] Gupta V, (2000) Legal Protection of Databases, *Malaysian J. of Library & Information Science*, **5** (2), 19-29.

[2] Agrawal R, Haas PJ, Kiernan J (2003) Watermarking relational data: framework, algorithms and analysis, *The VLDB J.*, , **12** (2), 157-169.

[3] Cox IJ, Miller ML, Bloom JA (2002) Digital Watermarking, *Morgan Kaufmann Publ., by Academic Press*.

[4] Chan C, Cheng L (2004) Hiding Data in Images by simple LSB substitution, *Pattern recognition*, **37**(3), 469-474. doi: 10.1016/j.patcog.2003.08.007.

[5] Lie W, Chang L (2006) Robust and high quality time domain audio watermarking based on low frequency amplitude modification, *IEEE Transactions on Multimedia*, **8**(1), 46-59. doi: 10.1109/TMM.2005.861292.

[6] Chan PW, Lyu MR (2003) A DWT –based Digital Video watermarking scheme with error correcting code, *In* proceedings of Inf. & comm. secur., LNCS, 2836, 202-213.

[7] Sion R, Atallah M, Prabhakar S (2004) Rights Protection for Relational Data, *IEEE Transactions on Knowl. & Data Eng.* **16**(12), 1509-1525.

[8] Directive 96/9/EC of the European Parliament and of the Council of 11 March 1996 on the legal protection of databases *Official Journal L 077, 27/03/1996 pp: 20-28.* Availableat:*http://eurlex.europa.eu/LexUriServ/LexUriServ.do?Uri=CELEX:31996L009: EN:HTML.*

[9] Black's Law Dictionary,(2009) 9th Edition, Database, ©thompson Reuters.

[10] United States Copyright Office, "Legal protection for databases", August 1997. Available at: *http://www.copyright.gov/reports/dbase.html.*

[11] Maggon H. Legal Protection of Databases: An Indian Perspective, 11, Journal of INTELL. PROP. RIGHTS, 140, pp: 140-144, 2006, *http://nopr.niscair.res.in/bitstream /123456789/3559/1/JIPR%2011(2)%20140-144.pdf,* http://papers.ssrn.com/sol3/papers. cfm?abstract_id =1398303.

[12] Saksena, Hailshree, Doctrine of Sweat of the Brow (May 3, 2009). Available at SSRN: http://ssrn.com/abstract=1398303 or *http://dx.doi.org/10.2139/ssrn.1398303.*

[13] Sinha S, Bench S, Sikri A. Eastern Book Company And Ors. vs D.B. Modak & Ors. & Mr. Navin J. on 27 September, 2002, *http://www.indiankanoon.org/doc/377266/.*

[14] Traditional Knowledge Digital Library, Collaborative project of CSIR & Dept. of AYUSH. Available at: *http://www.tkdl.res.in/tkdl/lang*default/common/Home.asp?

[15] WIPO News and Events, WIPO and India Partner to Protect Traditional Knowledge from Misappropriation, PR/2011/682, Geneva/Delhi, 22.03.2011, Available at:*http://www.wipo.int/pressroom/en/articles/2011/article_0008 .html*

[16] Indian Copyright act amendment (2012) www.wipo.int/edocs/lexdocs/laws/en/in/ in066een.pdf.

[17] WIPO, WIPO Copyright Treaty, Adopted in Geneva on Dec 20, 1996. Available at: *http://www.wipo.int/treaties/en/ip/wct/trtdocs_wo033.html*

[18] Digital Millennium Copyright Act, To amend title 17, United States Code, to implement the World Intellectual Property Organization Copyright Treaty and Performances and Phonograms Treaty, and for other purposes, Public Law 105-308, Oct-28, 1998.

[19] Australian Government ComLaw, Copyright amendment act-2006, C2006A00158, Available at: *http://www.comlaw.gov.au/Details/C2006A00158*

[20] Halder R, Pal S, Cortesi A (2010) Watermarking techniques for relational databases: Survey, classification and comparison, *J. of Universal Comp. Science*, **16**(21), 3164–3190.

[21] Khanduja V, Chakraverty S, Verma OP (2016) Ownership and Tamper detection of Relational Data: Framework, Techniques and Security Analysis, published as the chapter in the book titled: *Embodying Intelligence in Multimedia Data Hiding* at Science gate Publishing, pp:21-36, DOI: 10.15579/gcsr.vol5.ch2.

[22] Mehta BB, Aswar HD (2014) Watermarking for security in database: A review, *In* Proceedings of IEEE IT in Business, Industry & Government (CSIBIG), pp:1–6.

[23] Xie MR, Wu CC, Shen JJ, Hwang MS (2016) A Survey of Data Distortion Watermarking Relational Databases, *International Journal of Network Security*, vol. 18(6),1022-1033.

[24] Agrawal R, Kiernan J (2002) Watermarking Relational Databases, *In* proceedings of 28th VLDB conference, China, 155-166.

[25] Xinchun C, Xiaolin Q, Gang S (2007) A Weighted Algorithm for Watermarking Relational Databases, *Wuhan University J. of Natural Sciences*, **12**(1), 79–82.

[26] Farfoura ME, Horng SJ, Lai JL, Run RS, Chen RJ, Khan MK (2012) A blind reversible method for watermarking relational databases based on a time-stamping protocol, *Expert Systems with Applications*, **39**(3), 3185–3196.

[27] Zhou X, Huang M, Peng Z (2007) An additive-attack-proof watermarking mechanism for databases' copyrights protection using image, *In* Proceedings of ACM symposium on applied computing :254–258.

[28] Sun J, Cao Z, Hu Z (2008) Multiple watermarking relational databases using image, *In* Proceedings of IEEE Multimedia and Inf. Technology :373-376.

[29] Wang H, Cui X, Cao Z (2008) A Speech Based Algorithm for Watermarking Relational Databases, *In* Proceedings of Int. Symposium on Inf. Processing, 603–606.

[30] Shehab M, Bertino E, Ghafoor A (2008) Watermarking Relational Databases Using Optimization-Based Techniques, *IEEE Transactions on Knowl. & Data Eng.*, 20(1): 116-129.

[31] Khanduja V, Verma OP, Chakraverty S (2015) Watermarking Relational Databases using Bacterial Foraging Algorithm, *Multimedia tools & Applications*, Springer, 74(3): 813-839, DOI: 10.1007/s11042-013-1700-9.

[32] Iftikhar S, Kamran M, Anwar Z (2015) RRW-A robust and reversible watermarking technique for relational data, *IEEE Transactions on Knowledge and Data Engineering*, 27(4):1132–1145.

[33] Odeh A, Al-Haj A (2008) Watermarking Relational Database Systems, *In* Proceedings of the Applications of Digital Inf. & Web Tech. :270-274.

[34] Al-Haj Ali, Odeh A (2008) Robust and blind watermarking of relational database systems, *J. of Computer Science*, 4(12),1024-1029.

[35] Hanyurwimfura D, Liu Y, Liu Z (2010) Text format based relational database watermarking for non-numeric data, *In* Proceedings of IEEE ICCDA: 312-316.

[36] Sion R (2004) Proving ownership over categorical data, *In* Proceedings of ICDE: 584–595.

[37] Sion R, Atallah M, Prabhakar S (2005) Rights protection for categorical data, *IEEE Transactions on Knowl. & Data Eng.g*, **17**:912–926.

[38] Li Y, Guo H, Jajodia S (2004) Tamper detection and localization for categorical data using fragile watermarks, *In* Proceedings of ACM workshop on Digital Rights Management:73-82.

[39] Guo H, Li Y, Lui, Jajodia S(2006) Fragile watermarking scheme for detecting malicious modifications of database, *Elsevier, Information Sciences,* 176(10):1350–1378.

[40] Khan A, Husain SA (2013) A fragile zero watermarking scheme to detect and characterize malicious modifications in database relations, *The Scientific World J.,* Article ID 796726: 1-16.

[41] Khataeimaragheh H, Rashidi H (2010) A Novel Watermarking Scheme for Detecting and Recovering Distortions in Database Tables, *Int. J. of Database Management Systems,* 2(3): 1-11.

[42] Camara L, Li J, Li R, Xie W (2014) Distortion-Free Watermarking Approach for Relational Database Integrity Checking, *Hindawi Publishing Corporation*, *Mathematical Probl. in Eng.*:1-10 DOI: http://dx.doi.org/10.1155/2014/697165.

[43] Bhattacharya S, Cortesi A (2009) A distortion free watermark framework for relational databases, *In* Proceedings of Software & Data Technologies, 2: 229-234.

[44] Khanduja V, Chakraverty S, Verma OP (2016) Enabling Information Recovery with Ownership using Robust Multiple Watermarks, *J. of Inf. Secur. & Applications, Elsevier*, 29: 80-92. DOI: 10.1016/j.jisa.2016.03.005.

[45] Khanduja V, Chakraverty S, Verma OP, Singh N (2014) A Scheme for Robust Biometric Watermarking in Web Databases for Ownership Proof with Identification, *In* Proceedings of Active Media Technology, LNCS 8610, 212-215.

[46] Schneier B. (2008) Applied Cryptography, protocols, algorithms and source code in C, 2nd edn. Wiley-India.