

## A Big Data Security Method Using Modulus Operator

<sup>1</sup>R.A. Archana, <sup>2</sup>Ravindra S. Hegadi and <sup>3</sup>T.N. Manjunath

<sup>1</sup>R&D Centre, Bharathiar University, Coimbatore, Tamil Nadu, India

<sup>2</sup>School of Computational Sciences, Solapur University, Maharashtra, India

<sup>3</sup>Department of ISE, BMS Institute of Technology, Bangalore, Karnataka, India

---

**Abstract:** Due to internet of things and social media platforms, raw data is getting generated from systems around us in three 60° with respect to time, volume and type. Social networking is increasing rapidly to exploit business advertisements as business demands. In this regard, there are many challenges for data management service providers, security is one among them. Data management service providers need to ensure security for their privileged customers in providing accurate and valid data. Since, underlying transactional data have varying data characteristics such huge volume, variety and complexity, there is an essence of deploying such data sets on to the big data platforms which can handle structured, semi-structured and un-structured data sets. In this regard we propose a data masking technique for big data security. Data masking ensures proxy of original dataset with a different dataset which is not real but looks realistic. The given data set is masked using modulus operator and the concept of keys. Our experiment advocates enhanced modulus based data masking is better with respect to execution time and space utilization for larger data sets when compared to modulus based data masking. This research will help big data developers, quality analysts in the business domains and provides confidence for end-users in providing data security.

**Key words:** Big data, data security, data masking, huge volume, utilization, India

---

### INTRODUCTION

In today's business world, business organizations are more tend to use social media to exploit the business opportunities, social media applications are increasing day by day to make all their transactions and communications easy. Because of such transactions, huge volumes of heterogeneous data sets are generated. All these data sets need to be stored on big data platform-Hadoop for data analysis/predictive analysis. Business domains like banking, retail insurance, finance and security are using intelligent applications to make their business effectively. Information management solutions are no longer just a business enabler but an integral part of providing enhanced customer services. Information management is crucial for management of data efficiently and effectively. However, an overlooked risk for various domain in security space for data analysis. The various factors on security aspects such as types of real data used in software development, testing information security. Use of cloud computing, distributed computing, outsourced services and which experiences data breaches involving real consumer data. In this connection data masking provides data protection to

address privacy concerns. Data masking replaces sensitive data with a non-sensitive substitute but does, so in a way that preserves the integrity of the data. This means masked data can be used to facilitate business processes without changing the supporting applications, databases or data storage facilities which enables you to remove the risk without breaking your business. Securosis research has developed five laws for data masking: masked data should not be reversible, masked data should be representative of the original data set, masked data should maintain application and database integrity, non-sensitive data should be masked only if it can be used to re-create or tie back to sensitive data and data masking routines must be repeatable. One-off masking is both ineffective and impossible to maintain. Today's information technology environments are highly dynamic and masking routines need to keep pace. InfoSphere optim data privacy provides a comprehensive set of data masking techniques to support data privacy and compliance requirements. For the first time, you can mask data across platforms, across data sources using a standard and repeatable process to ensure data privacy without impacting the stability of your applications with greater ease and unparalleled scalability and performance.

With InfoSphere Optim Data Privacy, you mask and move. Masking and moving allows you to extract and mask data and then insert or load the data into one or more destinations. Masking in place allows you to de-identify data and replace existing values. InfoSphere optim data Privacy provides the most comprehensive set of data masking techniques on the market. The method you use will depend on the type of data you are masking and the result you want to achieve. Out-of-the-box capabilities for specific data types are included such as random or sequential number generation, string literal substitution, concatenating expressions, arithmetic expressions, lookup values and user-defined functions, to name a few. Some examples of situations in which masking techniques can be applied includes are data at rest or data in flight, relational data, flat files and data sets such as IBM IMS or VSAM, data being transformed through an Extract, Transform and Load (ETL) tool, data accessed in SQL queries inside a database, data in reports and documents, data inside applications, data moving to in and from big data platforms such as Hadoop, data used for testing big data environments, data used for analytics applications for example, PureData analytics or teradata and data used for testing data warehouses. Some of the benefits of data masking focus on data security and privacy to deliver significant value such as prevent data breaches, ensure data integrity, reduce cost of compliance and protect privacy.

**Literature review:** In contemporary information technology trend, social media is one among the top communication media to share the opinions of likes and dislikes. According to featured insights, global, media and entertainment report internet users spend more time with social media sites than any other type of site. At the same time, over the next 5 years, we project mobile broadband penetration to overtake fixed broadband, rising to 58.3% of the total in 2019, from 32.7% in 2014. Fixed broadband penetration growth will slow over the same period, rising by just 6.3% points, from 43.6% in 2014 to 49.9% in 2019 (Fig. 1).

In this regard, data security is very important for any business marketing across social media sites intruder or insider may steal the data which imbalances the business in the business market. The process of obscuring specific data elements within data stores is called as data masking. It ensures that sensitive data is replaced with realistic but not real data. The goal is that sensitive customer information is not available outside of the authorized environment. Data masking is an effective strategy in reducing the risk of data exposure from inside and August 16, 2018 outside of an organization and should be

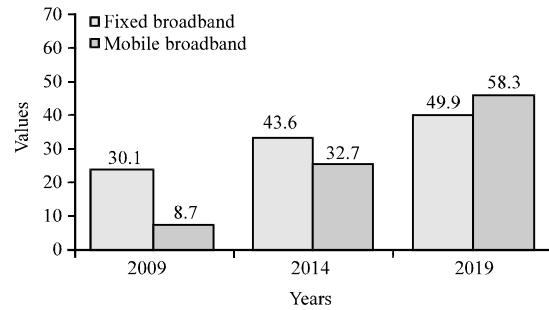


Fig. 1: Global broadband penetration (%); (Mckinsey and Company, Wilkofsky Gruan Associates)

considered a best practice for curing non-production databases (Manjunath *et al.*, 2011 a-c; Xiao-Bai and Luvai, 2009). The following researchers have exertion on data masking algorithms for various business domains in their perspectives. Muralidhar *et al.* (1995) proposed his research on “Random Data Perturbation for Non-normal Data” in 2000, Proceedings of Annual Meeting of the Decision Sciences Institute. Sarathy *et al* presented his research on “The Two Step Data Shuffle: A New Masking Procedure,” in Invited seminar presented to the Census Bureau and the Washington Statistical Society in 2002. Later in 2003 they gave the idea of “The Data Shuffle: A New Masking Procedure for Numerical Data,” in 8th INFORMS Computing Society Conference Sarathy, R. and K. Muralidhar gave the idea of “Data Masking-Problems, Solutions and Opportunities,” in TRDDC-TCS, Pune India in 2006 (Domingo-Ferrer and Mateo-Sanz, 2002; Bonifati *et al.*, 2001). Muralidhar and Sarathy (1999; 2008) gave the idea of “Privacy Violations in Accountability Data Released to the Public by State Educational Agencies” in Federal Committee on Statistical Methodology Research Conference in 2009. Ravikumar *et al.* (2011a, b) in (IJCSIS) International Journal of Computer Science and Information Security, Vol. 9, No. 8, August 2011, Data Masking Techniques with Random Replacement using data volume which is difficult for hackers to steal the data. Ricardo Jorge Santos *et al.* International Joint Conference of IEEE TrustCom, 2011, showed how balancing Security and Performance for Enhancing Data Privacy in Data Warehouses using modulus based method. Muralidhar and Sarathy (1999, 2008) states “Interval Responses for Queries on Confidential Attributes: A Security Evaluation.” Journal of Information Privacy and Security, 9(1), 3-16, 2013. A white paper by camouflage data masking specialist, titled “A Proactive Approach to Data Security for Cloud-Based Testing and Development”, May, 2014, emphasis any cloud-based application development offers organizations many tangible benefits,

yet, organizations struggle with how to work with data in the cloud-big data while complying with key regulations and meeting data security requirements. Data masking offer organizations a way to guard data for big data application development/testing using confirmed expertise while extenuating the risk of a security breach. Vishnu *et al.* (2014) international Journal of Computer Applications, recent Advances in Information Technology, 2014, proposed an effective data warehouse security framework which highlights on the usage of modulus operator in data security for data warehouse system (Vishnu *et al.*, 2014). No literatures found on creating uniform data security framework using enhance modulus based for big data. Hence, hereby a model is proposed which can be uniformly used across the industry for data security on big data environment which are business critical. According to recent Gartner research, unstructured data accounts for at least 80% of an organization's data. If left unmanaged, the sheer volume of unstructured data that's generated each year within a company can be costly in a number of ways, ranging from security vulnerabilities to compliance risks. The new era of computing has arrived: organizations are now able to process, analyze and derive maximum value from structured, unstructured and streaming data in real time. However, in the rush to achieve new insights are privacy concerns being neglected how can you support business goals while also, ensuring the privacy of sensitive data with the average cost of security-related incidents in the era of big data estimated to be over USD40 million, according to this Aberdeen Group Research Brief, you can't afford to ignore data privacy as a top requirement. With 2.5 quintillion bytes of data created every day, now is the time to understand sensitive data and establish business-driven privacy policies to keep customer, business, Personally Identifiable Information (PII) and other types of sensitive data safe. Remember, however, that different types of data will require different protection policies. For example, text, audio, log files and click streams have unique characteristics and challenges around privacy.

## MATERIALS AND METHODS

**Big data environment:** Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, curate, manage and process data within a tolerable elapsed time. Big data "size" is a constantly moving target as of 2016 ranging from a few dozen terabytes to many petabytes of data. Big data is not just about volume or rate of acquisition but also, heterogeneity/diversity, multiple levels of granularity,

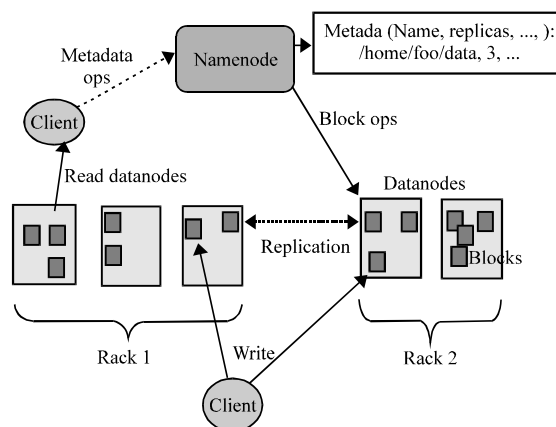


Fig. 2: HDFS architecture

multiple media and modalities, scientific disciplines, complexity, uncertainty incompleteness and representation opportunities. Big data presents unprecedented opportunities to accelerate scientific discovery and innovation, lead to new fields of inquiry that would not otherwise be possible, improve decision making, understand human and social processes, promote economic growth and improve health and quality of life. In this connection, we deploy variety of data on Hadoop platform for its usability based on the customer needs. Hadoop is an Apache top-level project being built and used by a global community of contributors and users. Hadoop was created by Doug Cutting and Mike Cafarella in 2005. Named in the remembrance of his son's toy elephant. The Apache Hadoop framework is composed of the following modules: Hadoop common which contains libraries and utilities needed by other Hadoop modules. Hadoop Distributed File System (HDFS) it is a distributed file-system that stores data on commodity machines, providing very high aggregate bandwidth across the cluster. Hadoop YARN, it is a resource-management platform responsible for managing compute resources in clusters and using them for scheduling of user's applications. Hadoop MapReduce, it is a programming model for large scale data processing. Figure 2 shows the HDFS architecture which has name node and corresponding data nodes in HDFS layer, similarly, MapReduce layer will have job tracker and task tracker.

**Mathematical model:** In our technique we make use of two keys where both are private, (i.e., the keys are known only to the authorized personnel). Consider a table say 'T' which has 'm' rows represented as (R1, R2, ..., Rm) and 'n' columns represented as (C1, C2, ..., Cn). Now, let us consider a value say (Ri, Cj) that has to be masked where 'Ri' and 'Cj' represents the row value and column value

respectively. Now the private keys used are K1 and K2. K1 is a random 128 bit random generated integer that remains constant for the whole table. K2 is another 128 bit random generated integer that remains constant for a single column whose value ranges between the maximum and minimum value of the column (in the above case for column Cj). Now for masking the data we use the formula:

$$(R_i, C_j)' = (R_i, C_j) - (K2 \% K1) + K2 \quad (1)$$

Now for de-masking we use the formula:

$$(R_i, C_j)' = (R_i, C_j) - (K2 \% K1) + K2 \quad (2)$$

Equation 1 is used for masking the data and Eq. 2 is used for de-masking.

**Proposed algorithm:** The proposed algorithm is used for data security using enhanced modulus based security with the concept of keys.

**Algorithm 1; Algorithm for data security using E-MOD:**

```

Begin algorithm
for each table n in the target databaseTD (1, ..., N)
fetch K1 private key common for the entire table (n)
for each attribute j in the table
    fetch K2, j private key common for entire attribute j
    for each item (Ri, Cj) in the tabel (n)
        fetch K3, i public key common for a tuple (I)
        find the length of the item (Ri, Cj), i.e., len = strlen ((Ri, Cj))
        If len NOT even
            append zero to (Ri, Cj)
        for each pair of character/digit in (Ri, Cj)
            Convert character into integer and store in numb
            apply masking formula for numb
            append to (Ri, Cj)' //every loop append with
                previous values
        end for]
    store (Ri, Cj)' in place of (Ri, Cj)
end for
end for
End algorithm
    
```

**Issues of data masking in security:** Business enterprises emphasized to covenant the following issues for a variety of data masking techniques such as: risk minimization: no matter what security measures are taken, there is always a degree of risk involved in handling a large amount of sensitive data (Ravikumar *et al.*, 2011). Data breaches can damage a company’s reputation increase liabilities and invite legal suits, accountability: data breaches create negative publicity, harm current and future business and damage organization’s reputation and the client’s confidence in it. It is crucial that the organization

stays accountable to all stakeholders, customers and employees and addresses their privacy needs effectively and regulatory norms: confidentiality and privacy norms demand the protection of data against theft. Compliance to all norms is essential to prove the organization commitment to its prestigious customers (Ravikumar *et al.*, 2011a, b). Data Masking Impediments. The points from a-e highlights the impediments of data masking methods in financial firms which are critical to business.

**Data utility:** Masked data should look and act like real data. Data must be fit for proper testing and development, application edits and data validation.

**Data relationships:** Must be maintained after masking on database level Referential Integrity (RI), application level RI, data integration (interrelated database RI).

**Existing business processes:** Must fit in with existing IT and refresh processes.

**Ease of use:** Must balance ease of use with need to intelligently mask data usable data that does not release sensitive information and knowledge of specialized IT/privacy topics and algorithmic importance should be pre-configured and built into the masking process.

**Customizable:** Solution/process must be capable of being tailored to specific needs of the clients (Dreibelbis *et al.*, 2008; Kimball and Caserta, 2004).

**RESULTS AND DISCUSSION**

The proposed methods provide flexibility around how the data will be masked and ensure that business rules of the enterprise application will not be impacted. After data segregation, the masking type will be decided based on the data such as substitution, replacement, multiplier, randomizer and shuffling, the same is illustrated below with example. Now the proposed technique E-MOBAT is compared with MOBAT taking the overheads-time taken for execution and storage space overheads. The results prove to be a lot in the favor of E-MOBAT proved to be better than MOBAT (Fig. 3, Table 1 and 2).

Thus, the above results prove that the E-MOBAT technique has a better performance curve than the existing MOBAT inturn making it better than the already existing algorithms such as AES128, 3-DES, etc.

Table 1: Storage space overheads

Size of data (MB)	MOBAT (MB)	E-MOBAT (MB)
2.50	3.015	2.19
5.00	6.223	4.30
10.0	12.250	8.22
20.0	24.420	16.50
40.0	48.440	32.90

Table 2: Execution time taken

Size of data (MB)	MOBAT (time sec)	E-MOBAT (time sec)
2.5	0.102	0.061
5.0	0.210	0.114
10.0	0.401	0.242
20.0	0.820	0.498
40.0	1.658	0.988

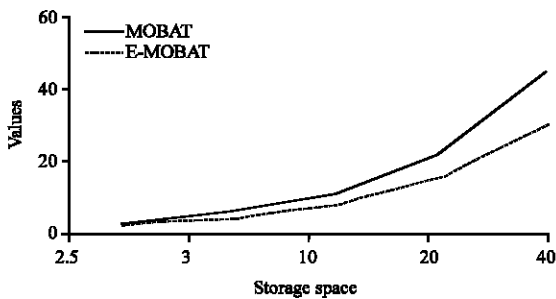


Fig. 3: Storage space overheads

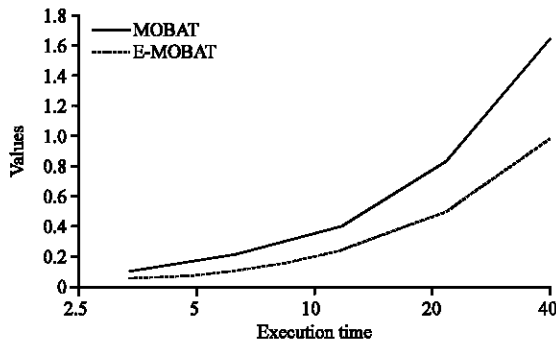


Fig. 4: Execution time overheads

Figure 4 shows the usage of primary memory (RAM) by both, the MOBAT and E-MOBAT to complete the execution of a query for varying number of values. The value in the x-axis represents the number of rows with each row having values from 9 different columns. The E-MOBAT uses the same amount of RAM as the MOBAT for smaller data-sets. But in real time scenarios involving millions of rows, the E-MOBAT functions at comparatively less RAM, than the MOBAT. The reduction in use of the primary memory is due to revision of the formulae used in the MOBAT. Proposed hybrid data masking method is a general approach that deals with the needs of security problems faced by various organizations when onsite-offshore business delivery model is used. Our hybrid data masking model framework

ensures two principles while operation is carried out: masking is not reversible. There is no way to reverse engineer the original data from the masked data and masked data is usable. For example when testing valid addresses the masked data must include valid zip codes not random numbers which fit the data type.

## CONCLUSION

In current information business market, masking offers unique value beyond other data security tools both in its ability to preserve complex data relationships while protecting data and its data management capabilities. Masking's combination of discovery, data set management, protection and control over data migration is unique. No other data security product provides all these benefits simultaneously. Masking reduces risk with minimal disruption to business systems. These characteristics suit masking to meeting compliance requirements. The rapid growth we have seen in the data masking segment spurred by compliance, risk and security demands has driven impressive innovate to capture increased customer demand.

## REFERENCES

- Bonifati, A., F. Cattaneo, S. Ceri, A. Fuggetta and S. Paraboschi, 2001. Designing data marts for data warehouses. ACM. Trans. Software Eng. Method., 10: 452-483.
- Domingo-Ferrer, J. and J.M. Mateo-Sanz, 2002. Practical data-oriented microaggregation for statistical disclosure control. IEEE. Trans. Knowl. Data Eng., 14: 189-201.
- Dreibelbis, A., H. Eberhard, M. Ivan, O. Martin and V.R. Paul *et al.*, 2008. Enterprise Master Data Management: An SOA Approach to Managing Core Information. Dorling Kindersley (India) Pvt. Ltd., Noida, India,.
- Kimball, R. and J. Caserta, 2004. The Data Warehouse ETL Toolkit: Practical Techniques for Extracting, Cleaning, Conforming and Delivering Data. John Wiley and Sons, New York, USA., ISBN-13: 9780764579233, Pages: 528.
- Manjunath, T.N., R.S. Hegadi and G.K. RaviKumar, 2011c. Design and analysis of DWH and BI in education domain. Intl. J. Comput. Sci., 8: 545-551.
- Manjunath, T.N., R.S. Hegadi and G.K. Ravikumar, 2011b. Analysis of data quality aspects in datawarehouse systems. Int. J. Comput. Sci. Inf. Technol., 2: 477-485.
- Manjunath, T.N., R.S. Hegadi and H.S. Mohan, 2011a. Automated data validation for data migration security. Intl. J. Comput. Appl., 30: 41-46.

- Muralidhar, K. and R. Sarathy, 1999. Security of random data perturbation methods. *ACM. Trans. Database Syst.*, 24: 487-493.
- Muralidhar, K. and R. Sarathy, 2008. A theoretical comparison of data masking techniques for numerical microdata. *Proceedings of the 3rd IAB Workshop on Confidentiality and Disclosure-SDC for Microdata*, November 20-21, 2008, Institute for Employment Research, Nuremberg, Germany, pp: 20-21.
- Muralidhar, K., D. Batra and P.J. Kirs, 1995. Accessibility, security and accuracy in statistical databases: The case for the multiplicative fixed data perturbation approach. *Manage. Sci.*, 41: 1549-1564.
- Ravikumar, G.K., J. Rabi, T.N. Manjunath, R.S. Hegadi and R.A. Archana, 2011b. Design of data masking architecture and analysis of data masking techniques for testing. *Int. J. Eng. Sci.*, 3: 5150-5159.
- Ravikumar, G.K., T.N. Manjunath, S.H. Ravindra and I.M. Umesh, 2011a. A survey on recent trends, process and development in data masking for testing. *Intl. J. Comput. Sci.*, 8: 535-544.
- Vishnu, B., T.N. Manjunath and C. Hamsa, 2014. An effective data warehouse security framework. *Intl. J. Comput. Appl.*, 2014: 33-37.
- Xiao-Bai, L. and M. Luvai, 2009. *Protecting Patient Privacy with Data Masking*. Firefly & Wisp Publishing, Homer City, Pennsylvania,.