

# Challenges for Big Data Security and Privacy

**B**

**M. Govindarajan**

*Annamalai University, India*

## INTRODUCTION

Big data refers to collections of data sets with sizes outside the ability of commonly used software tools such as database management tools or traditional data processing applications to capture, manage, and analyze within an acceptable elapsed time. Big data sizes are constantly increasing, ranging from a few dozen terabytes in 2012 to today many petabytes of data in a single data set. Big data creates tremendous opportunity for the world economy both in the field of national security and also in areas ranging from marketing and credit risk analysis to medical research and urban planning. The extraordinary benefits of big data are lessened by concerns over privacy and data protection.

As big data expands the sources of data it can use, the trust worthiness of each data source needs to be verified and techniques should be explored in order to identify maliciously inserted data. Information security is becoming a big data analytics problem where massive amount of data will be correlated, analyzed and mined for meaningful patterns.

Security of big data can be enhanced by using the techniques of authentication, authorization, encryption and audit trails. There is always a possibility of occurrence of security violations by unintended, unauthorized access or inappropriate access by privileged users.

To protect privacy, two common approaches used are the following. One is to restrict access to the data by adding certification or access control to the data entries so sensitive information is accessible to a limited group of users only. The other

approach is to anonymize data fields such that sensitive information cannot be pinpointed to an individual record. For the first approach, common challenges are to design secured certification or access control mechanisms, such that no sensitive information can be misconduct by unauthorized individuals. For data anonymization, the main objective is to inject randomness into the data to ensure a number of privacy goals (Xindong Wu et al., 2014).

## BACKGROUND

Today we are living in an era of digital world. With the rapid increase in digitization the amount of structured, semi structured and unstructured data being generated and stored is exploding. Usama Fayyad (2012) has presented amazing data numbers about internet usage like “every day 1 billion queries are there in Google, more than 250 million tweets are there in Twitter, more than 800 million updates are there in Face book, and more than 4 billion views are there in You tube”. Each day, 2.5 quintillion bytes of data are generated and 90 percent of the data in the world today were created within the past two years. The data produced nowadays is estimated in the order of zeta bytes, and it is growing around 40% every year. International Data Corporation (IDC) terms this as the “Digital Universe” and predicts that this digital universe is set to explode to an unimaginable 8 Zetabytes by the year 2015. The above examples demonstrate the rise of big data applications where data collection has grown tremendously and is beyond the ability of com-

monly used software tools to manage, capture, and process.

From a privacy and security perspective, the challenge is to ensure that data subjects (i.e., individuals) have sustainable control over their data, to prevent misuse and abuse by data controllers (i.e., big data holders and other third parties), while preserving data utility, i.e., the value of big data for knowledge/patterns discovery, innovation and economic growth.

Cloud protection alliance big data working group identify top protection and seclusion problems that need to confine for making the big data computing and infrastructure more secure. Most of these issues are linked to the big data storage and computation. There having some challenges which are related to secure data storage (Cloud Security Alliance White paper, 2012). Different security challenges related to data security and privacy are discussed in (A. A. Soofi et al., 2014) which include data breaches, data reliability, data accessibility and data support. Privacy is major concern in outsourced data. Recently, some controversies have revealed how some security agencies are using data generated by individuals for their own benefits without permission. Therefore, policies that cover all user privacy concerns should be developed. Furthermore, rule violators should be identified and user data should not be misused or leaked. The following sections describe some relevant challenges to security and privacy in the context of big data.

## **MAIN FOCUS**

### **Challenges for Big Data Security and Privacy**

With the proliferation of devices connected to the Internet and connected to each other, the volume of data collected, stored, and processed is increasing everyday, which also brings new challenges in terms of the information security. In fact, the currently used security mechanisms such as fire-

walls and DMZs cannot be used in the Big Data infrastructure because the security mechanisms should be stretched out of the perimeter of the organization's network to fulfill the user/data mobility requirements and the policies of BYOD (Bring Your Own Device). Considering these new scenarios, the pertinent question is what security and privacy policies and technologies are more adequate to fulfill the current top Big Data privacy and security demands (Cloud Security Alliance, 2013). These challenges may be organized into four Big Data aspects such as infrastructure security (e.g. secure distributed computations using MapReduce), data privacy (e.g. data mining that preserves privacy/granular access), data management (e.g. secure data provenance and storage) and, integrity and reactive security (e.g. real time monitoring of anomalies and attacks).

Considering Big Data there is a set of risk areas that need to be considered. These include the information lifecycle (provenance, ownership and classification of data), the data creation and collection process, and the lack of security procedures. Ultimately, the Big Data security objectives are no different from any other data types – to preserve its confidentiality, integrity and availability.

Being Big Data such an important and complex topic, it is almost natural that immense security and privacy challenges will arise (Michael & Miller, 2013; Tankard, 2012). Big Data has specific characteristics that affect information security: variety, volume, velocity, value, variability, and veracity. These challenges have a direct impact on the design of security solutions that are required to tackle all these characteristics and requirements (Demchenko, Ngo, Laat, Membrey, & Gordijenko, 2014). Currently, such out of the box security solution does not exist.

Cloud Secure Alliance (CSA), a non-profit organization with a mission to promote the use of best practices for providing security assurance within Cloud Computing, has created a Big Data Working Group that has focused on the major challenges to implement secure Big Data

6 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

[www.igi-global.com/chapter/challenges-for-big-data-security-and-privacy/183751](http://www.igi-global.com/chapter/challenges-for-big-data-security-and-privacy/183751)

## Related Content

---

### Detection of Shotgun Surgery and Message Chain Code Smells using Machine Learning Techniques

Thirupathi Guggulothuand Salman Abdul Moiz (2019). *International Journal of Rough Sets and Data Analysis* (pp. 34-50).

[www.irma-international.org/article/detection-of-shotgun-surgery-and-message-chain-code-smells-using-machine-learning-techniques/233596](http://www.irma-international.org/article/detection-of-shotgun-surgery-and-message-chain-code-smells-using-machine-learning-techniques/233596)

### An Efficient Source Selection Approach for Retrieving Electronic Health Records From Federated Clinical Repositories

Nidhi Guptaand Bharat Gupta (2022). *International Journal of Information Technologies and Systems Approach* (pp. 1-18).

[www.irma-international.org/article/an-efficient-source-selection-approach-for-retrieving-electronic-health-records-from-federated-clinical-repositories/307025](http://www.irma-international.org/article/an-efficient-source-selection-approach-for-retrieving-electronic-health-records-from-federated-clinical-repositories/307025)

### IS-Related Organizational Change and the Necessity of Techno-Organizational Co-Design(-In-Use): An Experience with Ethnomethodologically Oriented Ethnography

Chiara Bassetti (2012). *Phenomenology, Organizational Politics, and IT Design: The Social Study of Information Systems* (pp. 289-310).

[www.irma-international.org/chapter/related-organizational-change-necessity-techno/64689](http://www.irma-international.org/chapter/related-organizational-change-necessity-techno/64689)

### Mobile Technologies Impact on Economic Development in Sub-Saharan Africa

Adam Crossan, Nigel McKelveyand Kevin Curran (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 6216-6222).

[www.irma-international.org/chapter/mobile-technologies-impact-on-economic-development-in-sub-saharan-africa/184319](http://www.irma-international.org/chapter/mobile-technologies-impact-on-economic-development-in-sub-saharan-africa/184319)

### Do We Mean Information Systems or Systems of Information?

Frank Stowell (2008). *International Journal of Information Technologies and Systems Approach* (pp. 25-36).

[www.irma-international.org/article/mean-information-systems-systems-information/2531](http://www.irma-international.org/article/mean-information-systems-systems-information/2531)