# Big Data: Bottleneck Solution for Big Companies

[1] **Muhammad Raheel Zafar,** [2] **Muhammad Habib,** [3] **Kashif Razzaq,** [4]**Ahsan Raza Sattar**[4]

[1] Department of Computer Science, Lahore Garrison University, Lahore, Pakistan

[2] Department of Computer Science, Lahore Garrison University, Lahore, Pakistan

[3] Department of Computer Science, University of Agriculture, Faisalabad, Pakistan

[4] Department of Computer Science, University of Agriculture, Faisalabad, Pakistan

E-mail:  [1] raheel_zafar_iqbal@hotmail.com, [2]Ch.muhammadhabib@gmail.com , [3]mkashifrazzaq@gmail.com , [4] ahsan_raza@uaf.edu.pk

## ABSTRACT

Big data deals with enormous dissimilar types of data for the reason of data management with the help of business management analytical tools. It is a critical issue in current IT infrastructure that cannot be neglected in any future planning of business management. The structured and un-structured data are necessary to recognize and consolidate its nature with respect to data management. In this paper we have discussed the characteristics of data volumes, velocity, variety, value, legitimacy and complexity of the basic infrastructure of big data with respect to business analytics to understand its need for big companies to compete them with their competitors with efficient decision making and data handling within an appropriate time period.

**Keywords:**  *Big Data, Characteristics, Big Companies, Volume, Validity, Velocity, 4V.*

## 1. INTRODUCTION

Big data a buzzword deals with gigantic amount of structured data and always helpful in managing the unstructured data available in uncontrollable volume. With the passage of time it has increased in size and got deviation in traditional database infrastructure, management tools and techniques. Moreover the industries make a billion of transactions in every month within and outside its environment beside it demands wide storage space. With increasing data volumes the immediate information retrieval becomes tire some however quick retrieval is an important aspect for faster decision making. Generally, when stakeholder use the term big data means the requirement of technology to manage all kind of data and their storage as well. However, big data technology is utilized in search engines industry where customers demand high definition results in few seconds from vast distributed aggregation of tightly and loosely structured data [19].

Data management is a critical aspect for any organization. Any organization that covers the management flaws by the latest technology can lead the market. Approximately more than 1000 government and market leading companies in the world are adopting this technology and trying to maximize the utilization of their resources to increase their profit value. The development of new technologies and trends toward innovations has enhanced the value of big data. For instance, marketing teams analyzing the customer's taste by tracking their views, likes, comments or web clicks and enabling to improve the quality, campaign, stock and price of the products. Energy crisis regions such as developing countries, the industry or energy supply companies develop a sustainable energy saving and efficiency system through examining the utilization graph and reduce the energy consumption from unused sectors. Furthermore, the government and intelligence agencies as well as the search engines (Google) investigating the outgoing and incoming data and certify the authentication

of its validity. Similarly, geology institutes, metropolitan departments, mineral and oil, gas companies taking thousands numbers of outputs from sensor devices for generating results during its drillings and installation procedure to analyze for safety and efficiency. Generally, big data manages the large and complex data sets that are unrealistic to manage with traditional tools and techniques [25].

There are some general issues in the industrial sectors with respect to the increment of data amount on internet or cloud with rapid speed (Figure 1). These aspects are interrupting the way of success for big companies during the common transaction in between the basic industrial infrastructure i.e. decision making, management operation, stock handling, confidence level, finance, data manipulation etc. [11]
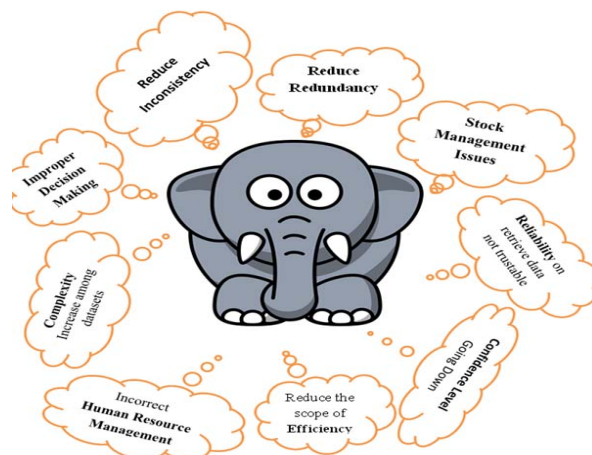


Fig. 1: Data sets increasing issues

Above figure elaborates the issues when the amount of the data sets increase. The aspects of the data relevant to the

quality disturb and cause of the poor management in any industrial sector. According to McKinsey & Company Research Institute survey report (2010), more than five billion phone users make trillion of transaction weekly and data is maintained by various telecom companies. Approximately, thirty billion contents are shared using Facebook every month. Data growth rate is above 40 percent through different projects including the 5 % data that is contributed by different IT companies through different projects. According to US Congress Library that stored 235 terabytes data from all over the world and 15/17 companies stored more data than US-Congress Library. With the passage of time, the size of the storage devices is shrinking and capable for storing enough data on it, 600 US dollar investment on storage device is enough for storing the audio music of the world.

Every industry enhances the scope of data importance for resources management and maintains the financial parameters. Data size increases rapidly in every sector and the judgment is made by utilizing the data and now it's a vital factor of the productivity improving of the industry. Big data constructs the value of data in several ways i.e. through creating transparency. Big data enables the researchers for experimentations to fulfill the organization needs, better performance and exposes variability as well as big data helps the higher level management to make decision by using some standardized algorithms instead of human recourses and it also open a new innovative way toward business, developing process and improving services [20].

In any organization, the main objective of the big data, that's why it is adopted specially with cloud computing technology, aims how to reduce the cost of data using big data technology. Time savings when data is stored in a manner able way. Developing the new products and services that generate results obey the base data that is collected from previous experiments. It also provides the 24/7 support for the online business and service provider, making internal organizations decisions. These are the aspects considered by the organization willing for implementing this technology in their business environment [7].

Higher level management is worried about the amount of data that is excessively increased. Beside it, what happened in worst cause when it create delays in decision with the growth of data among the departments and this gap hardly disturb the sale and purchase core modules of the system due to mismanage resources. If the human efforts are required for the validation of decision then why big data technology is implemented for such type of multidimensional data [4].

The above figure 2 elaborates that the data from different resources are gathered and then different technical and management teams utilize this data for reporting and decision making. Actually with respect to this scenario this environment is divided into three sections. Different kind of data is received and managed at middle layer where also have the support of different analytics tools for mature decision making and then store the data as well as results in next layer. Upper layer contain the stake holder, decided their role in the infrastructure.

In an organization, higher level management, stake holder, director, and technical staff concentrated on the big data technology and followed the following steps. Higher level authorities think about the big data implementation strategy when there are threats and data is uncontrollable. However, after analyzing all threat's factor, management analyze all their resources that organization or business required to handle this data and threats that occurs during handling this data. Once, an organization identifies their resources, data, and treatments, it's time to develop a strategic plan how utilize these resources to achieve our desired results and whose management tools and techniques are used during data management. At last, it is necessary to understand the complexity among the data and organization transactions to develop a more reliable, flexible and portable system [5].

Big data is an upcoming emerging technology whose popularity increases day by day in the world especially in huge business organization (Amazon), social media (Facebook, Twitter) and government sector. There data amount has not limit and increasing data bundles day by day with a rapid speed. Due to the huge amount of data it is much hard to maintain the quality of the data especially in the Asian industrial region, where quality gets less importance due to the high cost that an organization invests on quality rather than the production if compared it with the developed country. But in the big data, quality has vast worth because it is an expensive and sensitive infrastructure and compromising on the quality in this environment is not a good deal. Organizations feel hesitation to adopt due to the lack of security, availability and storage directly dependent on the quality.

In some aspect, it is considered that big data is the combination of social data and enterprise data as shown in in the below figure 3.



Fig. 3: Monetization of Data



Fig. 2: Data sets Analysis

The above figure elaborates the big data's stakeholder that how much their collaboration in data generation in different aspect and there accumulative result turned into big data. Structured and unstructured data is generated by stakeholder at different environment and store it on a platform where they get the access to the whole data according to their authentication and validation.

## 2. TYPE OF DATA

Big data is basically the collection of large amount of multi dimension structure and structure data under a manner able way for reliable and sophisticated decision making. With the passage of time the amount of data increases with a very uncontrollable emerging speed and approximately 300 quadrillion unstructured file store at different plate form whose amount is increased to double every year. Therefore organizations concentrate on the past, present and future aspect of the business; how to make their business more successful and ask few general question from there analyst that what happening, why is it happening. What happened, why did it happened what is likely and what should the organization do about it for improvement in their transactions. There are two types of data that majorly discuss in Big Data Technology;

### 2.1 Structured Data

Data stored in table in the form of fields and its attributes is known a structured data. Relational database and spreadsheets are the foremost example of structured in the era of modernization. Firstly, a data model that includes all type of business data, who's recorded the database, is built. This is followed by, how it was stored, its data type (numeric, bit, character, floating etc.) as well all the constraints that were implemented. The main advantages of the structured data [6] are;

- Storing of data at any platform easily inserted.
- Storing Strategy is easy to understand.
- Data analysis is much feasible for management.
- Accurate data retrieval is much an easier task
- Reduce complexity
- Time Saving



Fig. 4: Monetization of Data

In earlier days, due to the high cost, processor's limitation, storage, memory and processing constraints spreadsheets or relational database are used for maintaining the data in efficient and manner able way. Structured query language is normally used in industry for organizing the data; earlier develop by IBM in 1970 and then Oracle (Relational Software, Inc.) make some advancement according to the industrial requirements [13].

The text must be in English. Authors whose English language is not their own are certainly requested to have their manuscripts checked (or co-authored) by an English native speaker, for linguistic correctness before submission and in its final version, if changes had been made to the initial version. The submitted typeset scripts of each contribution must be in their final form and of good appearance because they will be printed directly. The document you are reading is written in the format that should be used in your paper.

### 2.2 Unstructured Data

It's totally a different concept with respect to structure in which data cannot be stored in traditional rows and columns. However, unstructured data file normally contain the audio, video, text files, word processing files, web pages, photos and business documents required for completing a transaction. It is considered unstructured as the data cannot be neatly inserted into database without any sorting. According to the experts, an organization is saving 80-90% of data files in unstructured environment and its amount increasing rapidly. Sometime its growth rate of storage is faster than structured data. Files of different formats are stored in distributed environment on hard disk; the major difference between structured and unstructured data [3].

## 3. RELATED WORK

The major focus of today modern science and business organization are big data and its analysis where the large number of data is generated through 24/7 transactions, audio video visual streaming, search engine, email, blogs, social media, RFID [31], camera, sensor, telecom packet and their applications that stored in database in structured and unstructured manner whose amount is growing up at tremendous rates, creating hurdles in insertion, deletion, updation, storage, analyzing, decision making and sharing through traditional database software tools[26]. The big data recognize as a buzz-word, also attract the industry about its glory in the global village when every organization wants centralized environment with the development of big data technology (continuous advancement technology) that increases the availability and production of human performance, also enhances the social activity of consumers. Cloud computing, green computing, grid computing, distributed computing, ubiquitous network and other current epoch technology provides an environment for the automation of the activity in data collection, visualizing, storing and processing [8].

The Decision making is an important phase in any industry through big data improve and this skill in any project management, increase efficiency, performance and effectiveness, however simply if institute implement the variety of different analysis tool and methodology to build an informative summaries data i.e. descriptive analysis generate standard reports, acknowledgement alerts and ad-hoc report;

predictive analysis is deal with statistical demonstration and forecasting by using it; and similarly prescriptive analysis aims to randomized and optimized the testing. Moreover it also helps the organization to perform some statistical analysis to recognize the organizational functionality [14].

The properties of big data can be recognized as "3V" (Volume, Velocity and Variety). At initial phase, huge amount of data bred from different resources i.e. Internet of Things (IoT) in which data retrieved over the internet from different sensor and network devices. Different trading organizations used millions RFID tag for shipment purpose of their inventories all over the world, similarly social media, a golden gate over a bundle of data stored daily and operated by millions of user i.e. Facebook, Twitter. First property of big data is volume refers to amount of data that is increased day by day with very high speed through different resources as discussed earlier. 2nd V is "variety" that deals with multiple types of data such as images, videos, audios etc. However in IoT, constant continuous streaming of structured and unstructured data is received with frequent speed that is generated from different devices. For instance, patients heart monitor, RFID, GPRS location information from phone that generated a lot of structured data but devices are not the only sources of data. Moreover, internet, social media, search engine, blogs are the way through miscellaneous set of structured and unstructured data obtained. An important property of big data is Velocity. It defines how much rapidly data move from one place to another over the network for both structured and unstructured data that required for manner able decision making. In the era of modernization as global has become a village and developed and IoT builds, every organization wish to capture the data frequently and make decisions for those "things" all around the world as soon as possible [24].

The world enter into the epoch of data driven where bundle of data acquired of different variety continuously moreover boost decision ability timely on the basis of available data from different platform is the key of successful business [18].The execution of any business process in any organization, especially where complex and large scale supply chains generate a high volume unstructured data, therefore immediately real time analysis is a difficult task. Integration of big data analysis with business process management architecture helps the industries analysis and improves the performance of business process [2].

The Data governance in big data is a good practice that cooperative to maintain a balance of data creation for high volume and exposure the risk factor especially for new born organization that helpful for unlocking the benefits and enhance value through application implementation of big data [27].The huge volume, complexity, data set productivity through multilevel, self-governing sources refers the big data. Due the rapid development of networks and data storage capacity, enhance the scope of the big data in all sciences, engineering, medical, nature and other several field [29].

## 4 Characteristics of Big Data

Now the days, the worth of big data should be identifying due to its following characteristics;

- Volume
- Velocity
- Value
- Veracity
- Varity
- Complexity

### 4.1 Volume

Here in big data, the term "Volume" identify the size of the data that is collected by the industry through different transaction, sensor devices etc. from previous large number of years and stored randomly on storage media (structured and unstructured Data). So, the amount of data increases rapidly in the era of computing and competition. However, when the data amount is growing up then the storage issues also occur specially for handling of unstructured data. With the improvement in software development methodology and techniques the hardware industry of storage devices is also growing up and the size of the devices is shrinking as well as cost also going down. These factors encourage the industrial sector to increase the data volume for analysis to make suitable decision [10].

Processing of very large amount of data for result is a complexity process. Principally, it is associated with different sensor devices, or forums for gathering the Digital data from different Libraries or resources whose growth rate is increasing day by day with unexceptional speed and potential used for further operations [28]. The term "Big data" itself define the volume and in the current era, data on hand is in petabytes and its amount expected in future to zetabytes. Amazingly, social media platform generate terabyte data every day and its amount is accidently increase with rapid speed that is hard to handle with the traditional methodology [17]. Bundle of data in an organization for building information is known as that is essential body as well as the assets that can be utilize for long period of time. However, with the increasing rate of data records, data value will also be diminishing [15].

The analysis of data growth rate identifies the increase 44% every year from 0.8zb to 35 zb till 2013-2020. Daily 15 terabyte data generated at twitter expected its going up to 70 tb per day till 2020. The below figure mention the accidental enhance of data on cloud that aspect analyze by IBM, according to them data from different organization increase and this increment ration obey the Moree's Law with respect to data increment ratio at different scenarios or platforms elucidated through this figure 5:
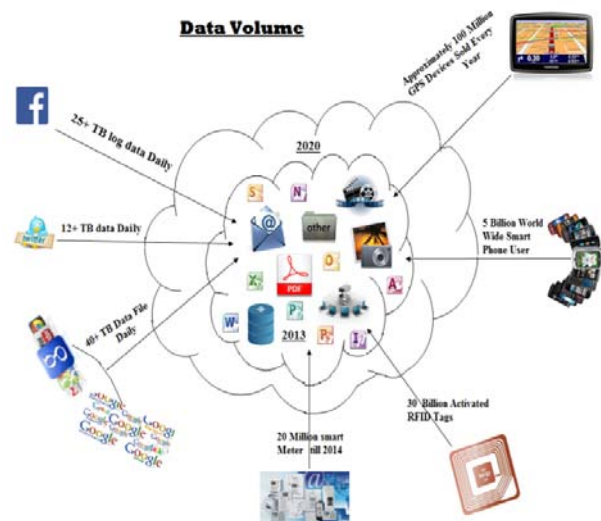


Fig. 5: Data Volume increasing Ratio

## 4.2 Data Velocity

In Big Data, the term "Velocity" means the speed of the data storage that received from different resources as well as it also deals with the data flow speed in between different platform. Consistency is maintain in big data because it deals with the incoming data that constantly store on storage device. For instance data generate from different sensor devices(Security Camera, GPS, Active RFID etc.) continuously store on storage medium. However, the traditional storage management techniques cannot achieve the desired performance when the data store continuously with a constant speed [17].

About twenty to thirty million person visit search engine platform like Google or Bing and search there result from all over the world immediately with in few second, here each and every user's click is recorded. In era of e-business, shopping is made easily, all products available at home. eBay or Amazon make the transaction procedure much faster than traditional purchase system as well also captured the items clicks and maintain logged of their user and allow them to make some immediate solution or decision. At different social media networks recorded its user each and every activity and behavior make some necessary for him according to its behavior, likes and dislikes. Cell phones, text message (email etc.), GPS, sensor data and RFID devices send storage signal continuously and big data provide a mechanism to save such packet with high speed without any interruption [9]. The speed of data writing on any storage medium is more important than the amount of data or volume of data because consistent and efficient speed is required for storing the continuous data received from different instrument or platforms etc. if storage occur with the high speed then then it make more sharp the organization against the competitors and protect the industry from any uncertain activity [21].

Velocity generally measures the speed it has great impact with the incremental issue of data how it handle. Velocity is also define that how frequently the data is generated (i.e. every Nano second, second, minutes, hour day week, month, years) by different resources and the processing speed also fluctuate among the user with respect to its requirements specification i.e. some kind of data need consistent frequency for processing real time data and some time it make necessary action when needed. However, divide these actions into three categories according to their storage nature: real time, occasional, frequent [30].

Big data and its analytics for real time environment containing the enormous focus toward data manipulation especially for telecommunication and social media and HP analyst elaborated the velocity aspect with reference to data volume. [16].

## 4.3 Varity

In the epoch of data revolution, with the amount of data volume incremental the variety of the data from different sources also increase that is uncontrollable to manage just like data volume in business to business transactions. A vast kind of data such as log files, images, videos, texts, emails, new domains, audio, video data, applications software etc. [23]. Big data variety variations at different platforms after every 60 seconds is shown in the below figure 6.
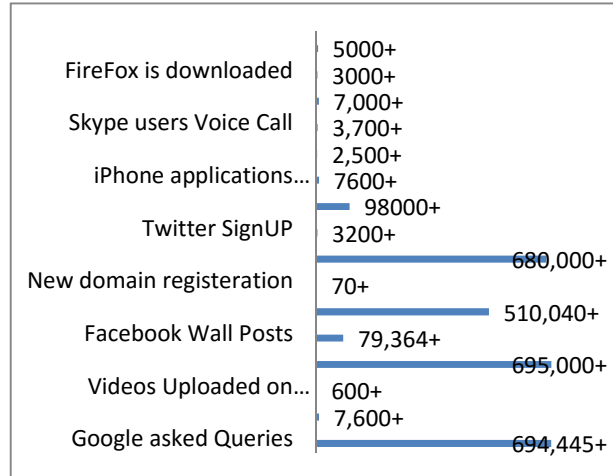


Fig. 6: Data Volume

The above figure elaborates that every day in every second, how much frequently the data of different variety is collaborated from different networks. People asked approximately 694,445 queries from Google and it helped them through analysis of different relevant platform. Now a days, Facebook is the most wanted, most precious social network, that handles more than 695,000 people in every minute updating their status, 510,140 commenting on them and furthermore 79,364 people make wall post on himself and there friend at Facebook in each and every minute. Similarly, approximately, 3200 new user signups at Tweeter and 98000 tweets after every 60 second and some more calculation seen in the above figure such as first two entities (WordPress and Firefox) shows that how much frequently the number of software developer and web user increasing in every minute [12].

Variety is the 3rd core characteristic of big data that represent different type of data. Storage and data analysis is varies from platform to platform like location coordinator, data sent from browser, simulators and video files etc. The aim of this characteristic is to sort all kind of data in a reliable manner that make it readable for user and help them to avoid from ambiguous results. Data received in big data is in unsorted form so a mechanism way is to be adopted t sort the data in a meaningful manner [1].

## 4.4 Data Value

Data value means how much data is usefulness for an enterprise for decision making and categorize them according to its usage. The aim of these measurements is to compute the insights, not the data file count. In data science, investigation technique is adopted that is useful to getting knowledge about data. On the other hand analytics sciences cover the predictive power of big data [15].

## 4.5 Data Veracity

Consistency in between such type large data is necessary to creating good impression in Big Data Technique. It is a fifth characteristic of big data. It deals with the uncertainty that happened due to data inconsistence and incompleteness, latency, ambiguity, approximation and deception etc. Suppose, entity 'A' send a message toward entity B and it

received an exact content that send by A in a mean time. If data loss from any Geo Location that there is no issue because it can be cover through other different ways. Hadoop, open source software that is used to handle large amount of data in a suitable time [1].

### 4.6 Complexity

Complexity identify the interconnectivity and independency in big data structure, similarly it held when a small change (or group of small changes) in few entity that occur a vast changes that ripple across or cascade through the system and strongly upset its behavior, or without any change at all [15].

All the above characteristics of big data discussion have their own importance according to their nature however, first four V's (4 V) are acting as core characteristics of big data that can be differentiated through the figure 7 that shown below.



Fig. 7: Characteristics of big data

In Figure 7 [22], try to develop the understanding about all core V's of big data through pictorial representation and their impact on data transactions that how these characteristics play their role. The above figures purify the concept among all the properties of big data and reduce the ambiguity in between of them. Volume only deals with the amount of data as well as velocity only concentrate with the speed of data transfer from one destination to another as soon as possible. On the other hand, in the scenario of data generation, a lot kind of data (Structured and un-structured) is received i.e. Text, Picture, Video, Mail and PDF etc. Veracity woks with the accuracy of the fetch data without any ambiguity, doubt occur due to any data in consistency and incompleteness.

## 5. Conclusion

Results are the main concentration of any companies that must be accurate, authentic, valid and concurrent. Big data implementation requires the basic understanding of it properties that depend on the architectural infrastructure capability of this cloud base atmosphere. Volume velocity, veracity and Varity are the core observation of big data.

## 6. Recommendations

Validation and verification of data is an important phase testing in software engineering that increases the authenticity of performed operations. It is a considerable aspect to validate the uploaded data and build a scenario to merge the same kind/ type of data (file, picture, video etc.) received from different user to manage the data volume in the big data. We try to identify the requirement elicitation phase with respect validation and verification of data in big data.

## REFERENCES

[1] A. Adrian, "Big Data Challenges", Journal of Database Systems, Vol. 4, No. 3, 2013, pp. 31-40.

[2] A. V. Baquero, and R. C. Palacios, "Business Process Analytics Using Big Data Approach", Journal of IT Professional, Vol. 15, No. 6, 2013, pp. 29-35.

[3] Blumberg, R., and S. Atre, 2003. The problem with unstructured data, DM REVIEW, 13 (1): 42-49.

[4] B. Brown, M. Chui, and J. Manyika, "Are you ready for the era of 'big data'?", Mckinsey Global institute, McKinsey & Company, 2011, pp: 1-12.

[5] J. Bughin, J. Livingston, and S. Marwaha, "Seizing the potential of 'big data', business technology office", International Institute for Analytics, McKinsey & Company, 2011, pp: 1-8.

[6] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber, "Bigtable: A distributed storage system for structured data", ACM Transactions on Computer Systems (TOCS), Vol. 26, No.2, 2008, pp. 4.

[7] T. H. Davenport, and J. Dyché, "Big Data in Big Companies", International Institute for Analytics, 2013, pp. 1-31.

[8] Y. Demchenko, P. Grosso, C. D. Laat, and P. Membrey, Addressing Big Data Issues in Scientific Data Infrastructure, International Conference on Collaboration Technologies and Systems (CTS), San Diego, Canada, 2013, pp: 48-55.

[9] L. Einav, and J. Levin, "The Data Revolution and Economic Analysis", in Conference on Innovation Policy and the Economy (No. w19035), National Bureau of Economic Research (NBER), 2013, pp. 1-29.

[10] A. Gupta, N. Mishra, J. Agarwal, and R. Patel, "A Study on Big Data", in International Conference on Cloud, Big Data and Trust, 2013, pp. 218-221.

[11] S. E. Hampton, C. A. Strasser, J. J. Tewksbury, W. K. Gram, A. E. Budden, A. L. Batcheller, C. S. Duke, and J. H. Porter, "Big data and the future of ecology", Frontiers in Ecology and the Environment, Vol. 11. No. 3, 2013, pp. 156-162.

[12] M. Heisterman, "In 60 Seconds [infographic]", 2014, http://www.astonishdesign.com/blog/60-seconds-infographic.

[13] S. A. Holcombe, R. H. Kohn, J. Knott, and R. Daniels, "Web-based royalty system and user interface", U.S. Patent (8,712,825), 2014.

[14] R. C. Joseph, and N. A. Johnson, "Big Data and Transformational Government", Journal of IT Professional, IEEE Computer Society, Vol. 15, No. 6, 2013, pp. 43 – 48.

[15] S. Kaisler, F. Armour, J. A. Espinosa, and W. Money, "Big Data: Issues and Challenges Moving Forward", in 46th Hawaii International Conference on System Sciences, IEEE Computer Society, 2013, pp. 995-1004.

[16] R. Kalkota, "Sizing "Mobile + Social" Big Data Stats", 2012,http://practicalanalytics.wordpress.com/2012/10/22/sizing-mobile-social-big-data-stats/

[17] A. Katal, M. Wazid and R. H. Goudar, "Big Data: Issues, Challenges, Tools and Good Practices", in 6th International Conference on Contemporary Computing (IC3), 2013, pp; 404-409.

[18] D. Keim, H. Qu., and K. L. Ma, Big-Data Visualization, Journal of IEEE Computer Graphics and Applications, Vol. 33, No.4, 2013, pp.20-21.

[19] P. Kumar, and K. Pandey, "Big Data and Distributed Data Mining: An Example of Future Networks", International

Journal of Advance Research and Innovation, Vol. 1, No.2, 2013, pp. 36-39.

[20] J. Manyika, M. Chui, B. Brown, J. Bughin, R. Dobbs, C. Roxburgh, and A. H. Byers, "Big data: The next frontier for innovation, competition and productivity", McKinsey GlobalInstitute(http://www.mckinsey.com/insights/mgi/research/technologyand_innovation/big_data_the_next_frontier_for_innovation), 2011, pp. 1-156.

[21] A. McAfee, and E. Brynjolfsson, "Big data: the management revolution", Harvard Business review, Vol. 90, No.10, 2012, pp. 60-68.

[22] T. Mole, How should we define big data? in Musings from a Big Data conference, 2013, http://www.intergen.co.nz/blog/Tim-Mole/dates/2013/8/musings-from-a-big-data-conference/

[23] B. Nedelcu, "About Big Data and its Challenges and Benefits in Manufacturing", Journal of Database Systems, Vol. 4, No. 3, 2013, 10-19.

[24] D. E. O'Leary, "Artificial Intelligence and Big data", Journal of IEEE Intelligent Systems, IEEE Computer Society, Vol. 28, No. 2, 2013, pp. 96-99.

[25] T. Rabl, S. G. Villamor, M. Sadoghi, V. M. Mulero, H. A. Jacobsen, and S. Mankovskii, "Solving big data challenges for enterprise application performance management", Proceedings of the VLDB Endowment, Vol. 5, No. 12, 2012, pp. 1724-1735.

[26] S. Sagiroglu, and D. Sinanc, "Big Data: A Review", in International Conference on Collaboration Technologies and Systems (CTS), San Diego, Canada, 2013, pp. 42-47.

[27] P. P. Tallon, Corporate Governance of Big Data: Perspectives on Value, Risk, and Cost, Journal of Computer, IEEE Computer Society, Vol. 46 No. 6, 2013, pp. 32-38.

[28] J. Wieczorkowski, and P. Polak, "Big data: Three-aspect approach", Online Journal of Applied Knowledge Management, Vol. 2, No. 2, 2014, pp: 182-196.

[29] X. Wu, X. Zhu, G. Q. Wu, and W. Ding, "Data Mining with Big Data" , IEEE Journal of Transactions on Knowledge and Data Engineering" , Vol. 26, No. 1, 2014, pp: 97-107.

[30] A. Zaslavsky, C. Perera, and D. Georgakopoulos, "Sensing as a service and big data.", International Conference on Advances in Cloud Computing (ACC), Bangalore, India, arXiv preprint arXiv:1301.0159, 2013, pp: 1-19.

[31] M. Habib, M. R. Zafar, S. Javed, and S. Ara, "A Sector Analysis for RFID Implantation: Technical Analysis for Integrated Security Enhancement Techniques",

International Journal of Engineering and Advanced Technology (IJEAT), Vol. 4, No. 1, 2014, pp. 132-136.

## AUTHOR PROFILES

**Muhammad Raheel Zafar**
MS (CS) 2014
M.Sc (CS) 2012
Research Assistant (IT) 2012-2014, Univ. Agri, FSD
Lecturer, Lahore Garrison University LHR Current.
i- Muhammad Habib, **M. Raheel Zafar**, Saima Javed, Shafaq Ara. A Sector Analysis for RFID Implantation: Technical Analysis for Integrated Security Enhancement Techniques
ii- Secure DNS from amplification attack by using. Modified Bloom Filters. Uzma Sattar.Talha Naqash. **M. Raheel Zafar**, Kashif Razzaq

**Muhammad Habib**
MS (CS) in progress
M.Sc (CS) 2012
System Administrator (IT) 2012-2014, WASA, FSD
Lecturer, Lahore Garrison University LHR Current.
i- **Muhammad Habib**, M. Raheel Zafar, Saima Javed, Shafaq Ara. A Sector Analysis for RFID Implantation: Technical Analysis for Integrated Security Enhancement Techniques.

**Kashif Razzaq**
M.Sc (CS) 2009
B.Sc (CS) 2005
Software Engineer 2014, CS Dept. Univ. Agri. FSD.
Secure DNS from amplification attack by using. Modified Bloom Filters. Uzma Sattar.Talha Naqash. M. Raheel Zafar, **Kashif Razzaq**

**Ahsan Raza Sattar**
Asst. Professor
Department of CS, Univ. Agri. FSD.