# Parallel Data Mining Architecture for Big Data

**Ms. Deepa P. Vaidya**                                   **Dr. Shrinivas P. Deshpande**

*Abstract* — **Big Data is concerned with the large-volume, complex, growing data sets with multiple and autonomous sources. With the fast development of networking, data storage, and the data collection capacity, Big Data is now rapidly expanding in all domains of engineering and science, including also the physical, biological and biomedical sciences. We need new tools and new algorithm to deal for this huge amount of data. The advances in hardware and software technologies have enabled us to collect, store and distribute large quantities of data on a very large scale. The process of discovering and extracting hidden knowledge in the form of patterns from these large data volumes is known as data mining. Data mining technology is not only a part of business intelligence, but is also used in many other application areas such as research, marketing and financial analytics.**

**Mining Big data has opened many new challenges and opportunities. Existing data mining techniques face great difficulties when they are required to handle the unprecedented heterogeneity, volume, variety, speed, privacy, accuracy and trust coming along with big data and big data mining. Extracting knowledge in the form of patterns from these massive growing data volumes in big data imposes a number of computational challenges in terms of processing time, memory, bandwidth and power consumption. These challenges have led to the development of parallel and distributed data analysis approaches. Various architectures have been suggested by the researchers, which mainly focus on mining the big data when it is stored in data repositories. This paper gives a conceptual model of parallel data mining architecture for big data. The architecture suggested shall process the big data or data stream using parallel data mining scheme prior to storage in data repositories.**

*Key Words* — **Big Data, Data Mining, Parallel Architecture.**

## I. INTRODUCTION

The Big Data is available from the heterogeneous, autonomous sources, in extreme large amount, which gets updated in fractions of seconds. [1] Therefore we can say that the Big data handles the real time data. Big data not only stores the petabytes or hexabytes of data in a data warehouse but it also provides the ability to formulate better decisions and take meaningful actions at right time.

The properties of big data are Volume, Velocity, Variety, Value and Varacity. Big Data requires processing high **Volumes** of data, which has unknown value, such as twitter or facebook data feeds, clicks on a web page, network traffic, sensor-enabled equipment capturing data at the speed of light, and many more. **Velocity** deals with a fast rate with which the data is received, stored and then acted upon for analyzing. The highest velocity data generally streams directly into the memory versus writing it to disk. Mobile application experiences huge user populations, increased network traffic, and the expectation for immediate

response. **Variety** includes the unstructured data types and semi-structured data types, such as text, audio, and video. This requires supporting metadata and additional processing to derive meaning. Unstructured data has same requirements as structured data, such as summarization, lineage, auditing and privacy. **Value** specifies the discovery of intrinsic value of data. There are a range of quantitative and investigative techniques to derive value from data – from discovering a consumer preference or sentiment, to making a relevant offer by location, or for identifying a piece of equipment that is about to fail. [2] **Veracity** refers to the inaccuracy of data like noise and abnormality. In deciding the big data strategy, it is necessary to carry out the data cleaning and processes so that the inaccurate data is not accumulated in the system. [3]

Daily 100 terabytes of data is uploaded in Facebook. Wal-Mart handles about 1 million customer transactions every single hour. The massive use of internet today by everyone brings a torrent of social media updates, sensor data from devices and a voluminous of e-commerce data. This means that every industry or organization is flooded with data, which can be extremely valuable, if it can be used to retrieve important information. [3]

The general architecture of big data has been suggested by organizations like Microsoft, IBM, Oracle, SAP, etc. Microsoft's big data ecosystem architecture shown in Fig. 1 is data-centric and is comprised of four main components: Sources (Variety, Velocity, Volume), Transformation (includes collection, aggregation, matching, data mining), Infrastructure (data storage or database software, servers, storage, and networking) and Usage (results). Microsoft depicts the flow of big data and possible data transformations from collection to usage. [4]
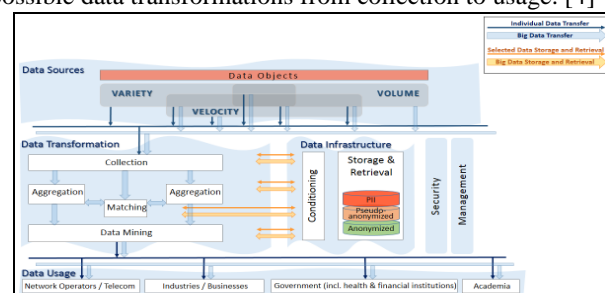


Fig. 1. Microsoft Big Data Ecosystem Reference Architecture

According to IBM the Big Data platform should support all of the data and must be able to run all the computations that are needed for its analytics. To achieve these objectives, any Big Data platform should address six key imperatives like Data Discovery and Exploration, Extreme Performance, Manage and Analyze Unstructured Data, Analyze Data in Real Time, Rich

Library of Analytical Functions and Tool Sets and Integrate and Govern All Data Sources. [4]

Oracle's Reference Architecture for Big Data gives a complete view of technical capabilities, how they fit together, and how they integrate into the larger information ecosystem. This helps to clarify Big Data strategy and to map specific products that support this strategy. Big Data reference architecture includes the infrastructure services, data sources, and information provisioning and analysis components. [4]

The Big Data Lifecycle Management (BDLM) model presents a new approach (see Fig. 2) to data management and processing in Big Data industry.
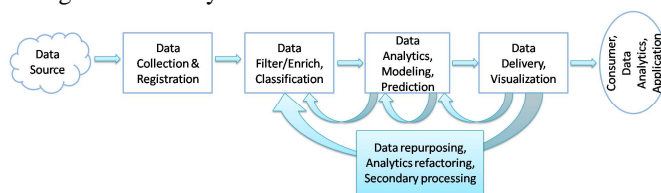


Fig. 2. Big Data Lifecycle Management

The BDLM needs the data storage and preservation at all the stages. It should permit data re-use/re-purposing and secondary research on the processed data and display the results. This can be achieved if the full data identification, cross-reference and linkage are implemented. During the lifecycle the data integrity, access control and accountability must be supported for the complete data. Data curation is an important component BDLM and must be done in a secure and trustworthy way. [5]

The ability to mine and analyze the big data from many sources, will give deeper and richer insights into business patters and trends. Thus it will help in driving the operational efficiencies and competitive advantage in manufacturing security marketing and IT. [6]

Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data owner. [7] It is a discipline, lying at the intersection of statistics, machine learning, data management and databases, pattern recognition, artificial intelligence, and other areas. [7,8,9,10]

The relationships and summaries derived through a data mining exercise are often referred to as models or patterns. Examples include linear equations, rules, clusters, graphs, tree structures, and recurrent patterns in time series. [7]

From the viewpoint of data mining, the big data has unlocked many new challenges and opportunities. The big data might be useful in various scenarios like detecting a fraudulent action while a person swipes credit card, or generating an offer to a customer when standing on a checkout line. Real Time Big Data Analytics for business personnel might be helpful for improving sales and for some scientists it can be a signal for an era in which machines begin to think and give humans like response using the vast amount of information which is being stored distributed.

Though big data bears greater value, it brings big challenges to extract the hidden knowledge and insights for knowledge discovery process. Current data mining techniques and algorithms are not ready to meet the new challenges of big data. The current data mining techniques when applied to big data shows limitations as they are centered on their inadequate scalability and parallelism.

In general, existing data mining techniques encounters many difficulties when they have to handle the unprecedented heterogeneity, volume, speed, privacy, accuracy, and trust coming along with big data and big data mining. By applying massive parallel architectures and novel distributed storage systems the existing techniques can be improved. Designing innovative mining techniques based on new frameworks/platforms with the ability to successfully overcome the aforementioned challenges will change and reshape the future of the data mining technology. [11]

In this paper a parallel data mining architecture is proposed for big data. The general architectures defined deals with the big data stored in data repositories. After storage the Data Mining is performed and models, rules and patterns are generated. But in case of Big data continuously data arrives at repositories in huge amount like in case of sensor networks data, twitter, facebook, etc. The traditional approach therefore has limitations as the big data mining is performed on the data stored in data repositories and when another bulky data arrives, again the mining has to be performed to generate accurate information for decision making. Therefore if parallel data mining is performed on the big data at the time of its arrival, the models and patterns can be reconstructed instead of dealing with the data repositories. Parallel data mining dramatically reduces the response time for data intensive operations on large databases associated with decision support systems.

## II. PREVIOUS WORK DONE

Yuri Demchenko et. al. [5] suggested that Big Data technology usage is increasing in the field of science and in industry. The big data focuses on data centric architecture and operational models. In their paper they discussed nature and origin of Big Data from different scientific, industry and social activity domains. They described a term called the Big Data Ecosystem which compromises of basic information/semantic models, architecture components and operational models. The authors proposed an improved Big Data definition and the Big Data Architecture Framework (BDAF) to address the Big Data Ecosystem aspects like Big Data Infrastructure, Big Data Analytics, Data structures and models, Big Data Lifecycle Management, Big Data Security. The paper analyses the requirements and provides suggestions about how the various components considered by the authors can address the main challenges in Big Data technology. The paper intends to offer a consolidated view of the Big Data phenomena and related challenges to modern technologies.

Catalin BOJA et. al. [12] analyzed the problem of storing, processing and retrieving meaningful insights from petabytes of data. The authors have provided a survey on current distributed and parallel data processing technologies and proposed an architecture that can be used to solve the analyzed problem. The

purpose of author's research is to implement the cluster analysis model. In this paper more stress is given on distributed files systems and the ETL processes involved in the distributed environment. According to them, using a multi-layer architecture to acquire, transform, load and analyze the data, ensures that each layer can be used perfectly for its specific task. Regardless of the various layers involved, the user is exposed to a user friendly interface provided by the front end application server. Once the storing and processing of the data is done then the relationships between different types of data are found.

Pushpanjali et. al [13] describes that the aim of Big Data Mining is to extract the useful information from large data sets or data streams. They discussed an overlay based parallel architecture which does the processing by using the overlay network and distributed data management. This can be useful in achieving the high scalability and service availability. In case of network disruption the parallel data mining architecture will not provide proper services. They suggested that an overlay network construction scheme based on physical network structure and containing nodes location and distributed task allocation scheme can be helpful to overcome this issue.

Felicity George et. al. [14] discussed a parallel data mining architecture that has following features like: ❶Scalable, ❷Capable of handling massive data sets, ❸Faster on 1 processor also, if sufficient memory is provided, ❹Portable because most of the modules are system independent, ❺Extensible, ❻Simple to use. The mining done is highly efficient capable of scanning millions of rows of data per second. The mining process is divided in two components, viz; Compaq's Data Mining Server (DMS) and Syllogic's Data Mining Techniques / Methods of Parallelizing (DMT/MP). Compaq's DMS, a parallel server provides a set of data mining primitives. These primitives are utilized by Syllogic's DMT/MP; a data mining client, which implements the actual data mining algorithms. The parallel architecture and the primitives used to operate on the data and use of these primitives by mining algorithms are discussed. Performance figures are specified for both the primitives and the high level mining algorithms.

## III. EXISTING METHODOLOGY

The Big Data Architecture Framework (see Fig. 3) comprises of the 5 components as mentioned below:

1. **Data Models, Structures, Types**: It should support variety of data types produced by different data sources and need to be stored and processed, which define the Big Data infrastructure technologies and solutions can to some extent.

2. **Big Data Management Infrastructure and Services**: It should support the Big Data Lifecycle Management (BDLC), provenance, curation, and archiving. BDLC should support the major data transformations stages such like data collection and registration; data filtering, classification; data analysis, modeling, prediction; data

3. delivery, presentation, visualization; and can be completed with the customer data analytics and visualization.

4. **Big Data Analytics and Tools:** It addresses also infrastructure components and functionalities related to the required data transformation.

5. **Big Data Infrastructure (BDI):** It deals with storage, compute, network infrastructure, and also sensor network and target/actionable devices.

6. **Big Data Security:** It should perform the following functions like protecting data in-rest, in-move, ensuring trusted processing environments and reliable BDI operation, providing the fine grained access control and protect users personal information. [5]
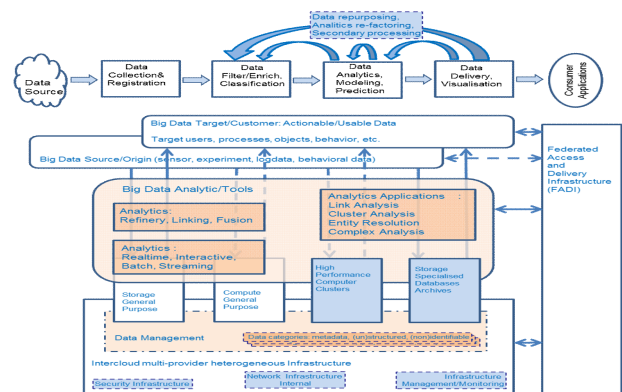


Fig. 3. General Big Data Infrastructure functional components

This architecture provides a view on the Big Data infrastructure for general data management, typically cloud based, and Big Data Analytics part which requires the high-performance computing clusters. [5]

The Distributed Parallel Architecture for Big Data, shown in Fig. 4, has following layers: [12]

- Input layer implements data acquisition processes by fetching data from different sources, reports, data repositories and archives. These sources use independent data schemes which is challenging problem. Bringing the data fetched in a common format is an intensive data processing stage. The architecture described has following layers:

- Data layer which implements distributed and parallel processing and responsible to store and process large datasets.

- User layer provides access to data and manage requests for analysis and reports.

- ETL intermediary layer between the first two main layers, collects data from the data crawler and harvester component, normalize data, discard unwanted information and bring it to a common format, requested by the parallel distributed DBMS that implements the Hadoop and MapReduce framework.[12]

The user will submit his inquiries using a front end application server, which converts them into queries and submit them to the parallel DBMS for processing. [12]
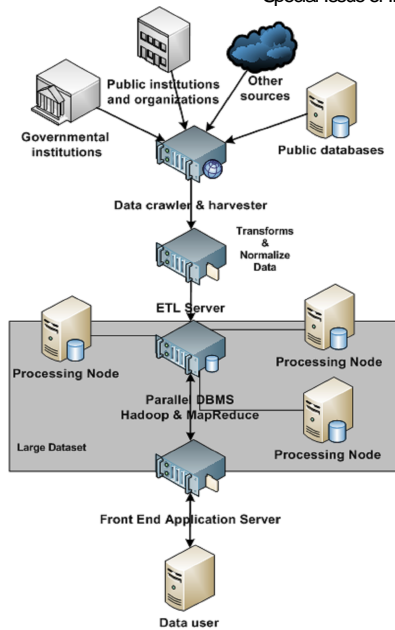
Fig. 4. Distributed Parallel Architecture for Big Data



Fig. 5. Mapping and Reduction processes in Overlay based parallel data mining architecture

The objective of the architecture is to separate layers, that does intensive data processing and link them to the ETL services. These ETL services will act as buffers or cache zones and does the transformation function. [12]

Overlay-based parallel data mining architectures improves the service availability against server break- downs. [13] All the servers in this architecture, executes management and processing functions. The overlay network is created using all the servers. In conventional architecture, processing nodes are similar to the master. This architecture provides continue servicing even if some nodes are removed from overlay network. An example of mapping and reduction processes in the architecture is shown in Fig. 4. The Mappers are randomly selected and found by the node when the request or message is brought (nodes B, C, and D in the Fig. 5). The mapper which finish the mapping process (node D in the Fig. 5) first becomes a reducer. It requests to other mappers to transmit the processed data towards itself. The request message can be forwarded by using flooding scheme. When the processed data are received from mappers, the reducer executes the reduction process and outputs the analyzed result. [13]
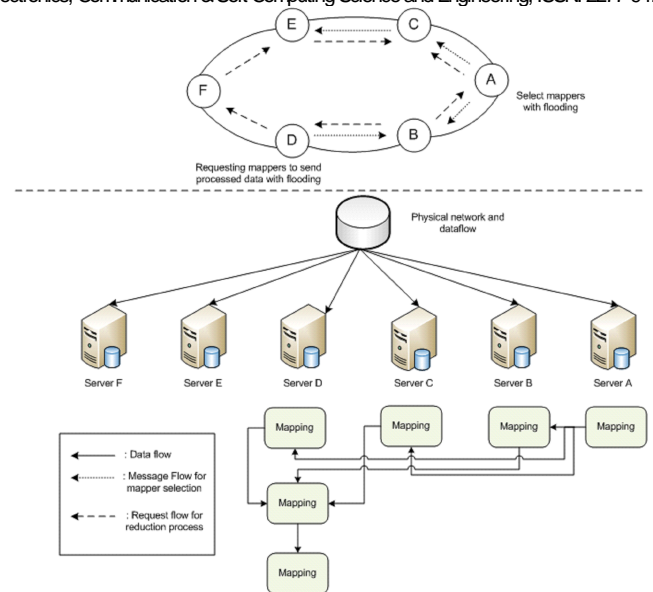
The overlay-based parallel data mining architecture is developed which is tolerant to physical network disruption, so that data mining is available at any time and at any location. [13]

A Parallel Data Mining Architecture for Massive Data Sets is described in another research paper. In Fig. 6 an overview of Data Mining Server is shown. The system consists of a manager and a number of servers. Each Server process a subset of the data. From the client data mining tool, request is received by the manager. If the request is simple query then it is a handled completely by the manager. If the request requires processing of the data, it is farmed out to the servers. The manager consolidates responses from the servers and returns the required results to the client. The manager also maintains a catalog of data detailing which DMS has access to for the current query. [14]
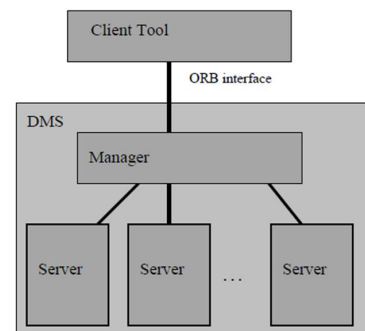


Fig. 6. Overview of DMS's Architecture

Fig. 7 shows the DMS components. Each server consists of four modules: [14]

- Factory: It Inputs data from variety of data sources into DMS. The factory converts the data into COREs (Column Represented Data Objects), DMS's internal data representation.
- Engine: It does all operations on DMS data, such as building

211

cross-tables and fetching rows of data to provide it to the front end client.

- Cache: It stores and manages COREs and other internal DMS objects.
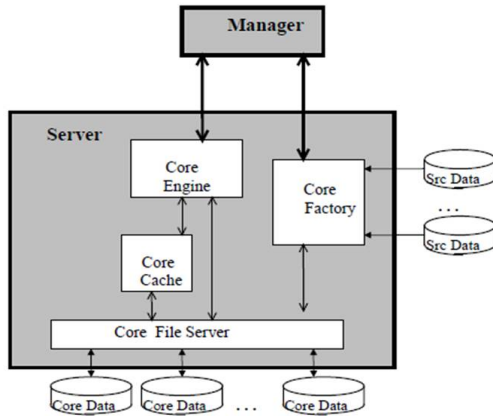- File Server: It manages the permanent storage of and access to COREs on disk. [14]



Fig. 7. Internal structure of DMS

A data-parallel record-based parallelization of data mining primitives like administrative, operational and data description/visualization is implemented by Data Mining Server (DMS). [14]

## IV. ANALYSIS AND DISCUSSION

The general Big Data Infrastructure (BDI) framework that has been suggested is based on various general architecture surveys. The general BDI capabilities, services and components included are: General Cloud based infrastructure, platform, services and applications, Big Data Management services and infrastructure, Registries, indexing/search, semantics, namespaces, Security infrastructure (access control, policy enforcement, confidentiality, trust, availability, privacy) and Collaborative environment (groups management). [5]

The Distributed Parallel Architecture for Big Data gives more emphasize over the distributed file systems and ETL processes done in distributed environment. The architecture includes layer like input, data, ETL and user. The objective of the architecture is to implement the cluster analysis model to acquire, store and process the large datasets. A combination of technologies like Hadoop-Map Reduce-Parallel DBMS solutions included in ETL layer can gain benefits from each technology. [12]

In the Overlay –based Parallel data mining architecture, the overlay network connectivity dramatically affects the service availability of data mining. There are various connectivity issues dealing works like context- cognizant, graph theory predicated, and intricate network theory predicated overlay network construction. These works make overlay networks, tolerant to minute-scale server breakdowns. They do not consider the sizable voluminous- scale server breakdowns, i.e., a physical network disruption. So this architecture is tolerant to network disruption and provides services of data mining without interruption. [13]

The parallel data mining architecture for massive data sets is simple and includes the components like client data mining tool, Manager, Data Mining Server (DMS) which process subset of data. This DMS architecture attains the data processing high speeds due to following factors:

- Effective data parallelization,
- Efficient data encoding into Column Represented Data Objects (COREs)
- Simple and optimized algorithms for COREs manipulation
- Data storage by column rather than by row, reducing disk access times and memory requirements for column based operations.
- Zoom-in functions enabling fast processing of subsets of the data, removing the need to scan the entire data. [14]

The analysis of continuous bombardment of data stream in big data and reconstructing the models and patterns using a parallel data mining architecture is not covered in any above architectures.

## IV. PROPOSED ARCHITECTURE

The social networking sites, sensor networks, online reservation systems, etc. generates continuously the structured and non structured data. The model and patterns generation from this continuous data stream in real time is a challenging task. In this paper a Parallel Data Mining Architecture for Big Data is proposed to handle such big data generated from data streams as shown in the Fig. 8.
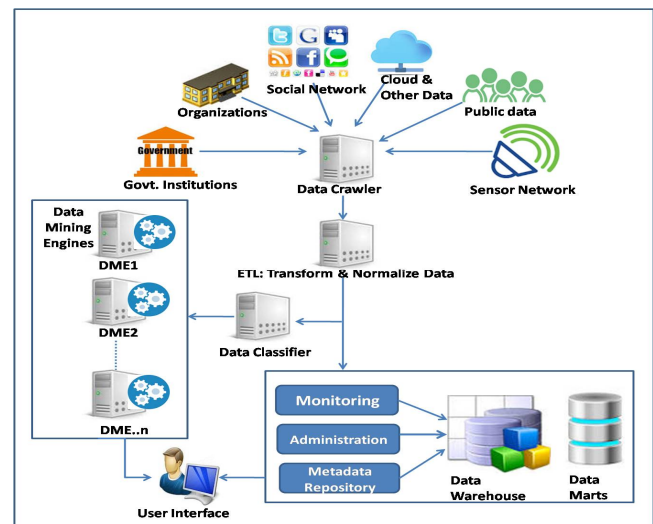


Fig. 8. Parallel Data Mining Architecture for Big Data

In this architecture the data is collected by Data crawler. Crawler creates and repopulates the data by downloading documents and files by navigating the web. The data gathered by Crawler can have various sources like Government institutions, Organizations, Social networking site, Public data, Sensor

Network, Cloud or other data sources. This collected data is processed by ETL server which will extract, transform, load, clean and refresh the data. This transformed data will be send to Data classifier and Data Warehouse (DW). The Data Warehouse will provide the data to user interface according to the query requested by the user. The Data Classifier will classify the data according to its types using various classification techniques as the data received from stream data sources is multimedia in nature. The Data Classifier will classify according to its types. It is required because the DMEs provided in this architecture performs Data Mining in parallel on different types of data. The models and patterns are generated by these parallel DMEs uses various Data mining algorithms. In general Data Mining approach, the models and patterns are generated from data repositories after storage using a training data set. But till the models and patterns are generated, the high volume, velocity and variety data stream keeps on coming in big data continuously. If we wait for storage of this data to Data Warehouse or data repositories and then apply data mining techniques for deriving the models and patterns, it will be time consuming and inaccurate. It desperately requires a real time modeling & mining integration with Data Warehouse organization. Therefore the Parallel Data Mining Engines layer will serve this problem. It will continuously analyze the data stream coming in Big Data and re-evaluate the models and patterns which will provide accurate results. According to the queries by the user interface the DME and DW will provide the results in the forms of report, graphs, etc. Thus this architecture will help in parallel data mining of Big Data.

## CONCLUSION

The volume of big data is increasing day by day. Big Data Technologies are targeting to process the high volume, velocity and variety data sets for extracting the meaningful data value and guarantees high-veracity of original data. Big data also concerns to obtain information that demand cost-effective, innovative forms of data analytics for improved insight, decision making, and processes control. Therefore big data has to be supported by new data models and new infrastructure services and tools that allow fetching and processing data in parallel from a variety of sources like social network, sensor networks, etc. and delivering data in a variety of forms to different data and information consumers and devices. Architectures that are developed and are defined to handle the big data, deals with the data repository stored in data sources. The architecture proposed in this paper will help in analyzing the data stream coming continuously in big data before storing in data repositories or data warehouse which will prove more efficient to generate accurate data models and patterns to derive meaningful results in less time.

## REFERENCES

[1] Rohit Pitre, Vijay Kolekar, "A Survey Paper on Data Mining With Big Data", (IJIRAE) International Journal of Innovative Research in Advanced Engineering, Vol. 1, Issue 1 (April 2014), pp. 178-180, ISSN: 2278-2311

[2] Oracle Enterprise Architecture White Paper, April 2015. Retrieved from http://www.oracle.com/technetwork/topics/entarch/articles/ oea-big-data-guide-1522052.pdf

[3] Kaushika Pal, Dr. Jatinderkumar R. Saini, "A Study of Current State of Work and Challenges in Mining Big Data", (IJANA) International Journal of Advanced Networking Applications, pp. 73-76, ISSN No. : 0975-0290

[4] Reference Architecture Subgroup, NIST Big Data Working Group (NBD-WG), Survey Report on Big Data Architecture Models, Ver. 1.2, (September 2013). Retrieved from http://www.google.co.in/url?q=http://bigdatawg.nist.gov/_uploadfiles/M0151_v3_3456779813.docx&sa=U&ei=K19UVbj2H86uuQSKkoHwCA&ved=0CB8QFjAB&usg=AFQjCNHlGMIPMChxRT3u5NvC6NZ7eDqPJA

[5] Yuri Demchenko, Canh Ngo, Peter Membrey, "Architecture Framework and Components for the Big Data Ecosystem", Ver. 0.2. Retrieved from: http://www.uazone.org/demch/worksinprogress/sne-2013-02-techreport-bdaf-draft02.pdf

[6] Intel White Paper on Mining Big Data in the Enterprise for Better Business Intelligence, (July 2012), Retrieved from http://www.computerworlduk.com/cmsdata/whitepapers/3526534/intel_big_data_-_mining_big_data_in_the_enterprise_for_better_business_intelligence.pdf

[7] Hand, D., Mannila, H. and Smyth, P. (2001), Principles of Data Mining [online], Retieved from: ftp://po.istu.ru/public/docs/other/_New/Books/Misc/Principles%20of%20Data%20Mining.pdf (Accessed 22 April 2012)

[8] Dunhum, M. (2006), Data Mining: Introductory and Advanced topics, New Delhi, India, Dorling Kindersley (India) Pvt. Ltd, pp. 4-275.

[9] Pujari, A. K. (2001), Data Mining Techniques, Hyderabad, India, University Press (India) Pvt. Ltd. pp. 42-149.

[10] Han, J. and Kamber, M. (2006) Data Mining Concepts and Techniques, Haryana, India, Elsevier Inc. pp. 1-685.

[11] Dunren Che, Mejdl Safran, and Zhiyong Peng, "From Big Data to Big Data Mining: Challenges, Issues and Opportunities". Retrieved from http://link.springer.com/chapter/10.1007/978-3-642-40270-8_1#page-1

[12] Catalin BOJA, Adrian POCOVNICU, Lorena B T GAN, "Distributed Parallel Architecture for "Big Data" ", Informatica Economic , Vol. 16, no. 2/2012,pp. 116-127, Retrieved From http://www.revistaie.ase.ro/content/62/12%20-%20Boja.pdf

[13] Pushpanjali, Jyothi S Nayak, "A Survey on Big Data, Data Mining and Overlay Based Parallel Data Mining", (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 6 (2) , 2015, pp. 1535-1538, ISSN:0975-9646

[14] Felicity George, Arno Knobbe, "A Parallel Data Mining Architecture for Massive Data Sets". Retrieved from http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.36.3420&rep=rep1&type=pdf

## AUTHOR'S PROFILE

**Ms. Deepa P. Vaidya**
Associate Professor & Assistant Co-ordinator,
P.G. Department of Computer Science & Technology,
D.C.P.E (Autonomous College), Shree H.V.P.Mandal Amravati, Maharashtra, India.
(e-mail: deepa_vaidya@rediffmail.com)

**Dr. S. P. Deshpande**
Associate Professor & Co-ordinator,
P.G. Department of Computer Science & Technology,
D.C.P.E (Autonomous College), Shree H.V.P.Mandal Amravati, Maharashtra, India
(e-mail: shrinivasdeshpande68@gmail.com)