

Distributed Computing in Big Data Frame Work : A Review

M. Bhargavikrishna

Research scholar
DeptofCSE

Prof. S. Jyothi

Dept of Computer science
Spmvv,Tirupati

Spmvv,Tirupatijyothi.spmvv@gmail.com
bhargavimandara@gmail.com

ABSTRACT

Big data analytics has attracted close attention from both industry and academic because of its great benefits in cost reduction and better decision making. As the fast growth of various global services, there is an increasing need for big data analytics across multiple data centers (DCs) located in different countries or regions. The data processing platform optimized for the geo-distributed computing environment. It supports a cross Dc data processing platform optimised geodistributed computing environment. To evaluate the performance the extensive simulations using real traces generated by a set of queries on Hive. The results show that proposal can reduce 55% interDC traffic compared with centralized processing by aggregating all data to a single data center. Public distributed computing is a type of distributed computing in which so called volunteers provide computing resources to projects. It introduces a distributed and selforganizing algorithm to build a management systems like health care. Research show that public distributed computing has the required potential and capabilities to handle big data mining tasks. It provides the foundation for future research required to bring back attention to this lowcost public distributed computing method and make it a suitable platform for big data and data mining analysis.

Keywords --

Big data analytics , geodistributed, Distributed algorithm, Self Organization, Data mining

INTRODUCTION

Many companies and organizations in today's world are interested in gathering data making tasks. It leads companies to capturing, storing and processing huge data sets, which in turn refers to a term called big data mining. Over time data sets become large and connected to many data points making data difficult to store and process. Big data mining is a computational process of discovering patterns in large data sets. It takes use of methods at the intersection of artificial intelligence, machine learning, statistics and database systems. Unfortunately, internal highperformance computing environments or similar traditional data management solutions are no longer capable of handling such amounts of data. Organizations do not usually have enough internal computational resources to satisfy the demand. More and more companies start to provide global services by deploying data centers (DCs) in different countries and regions. For example, Google runs its service across several geodistributed data centers connected by a dedicated WAN. Other companies, e.g., Netflix, deploy their services at Amazon's global cloud infrastructure EC2 that spreads across 11 regions over the world. These

e companies conduct big data analytics across the geodistributed computing and storage environment for risk evaluation, cost reduction, and new product creation. To deal with geodistributed big data analytics, several recent efforts have been made to create a virtual cluster across multiple DCs for big data processing. For example, Mandal et al have implemented and evaluated a Hadoop cluster across multiple clouds. Iridium has been proposed for low latency queries on geo-distributed big data.

BIG DATA

Big Data is a collection of data that is huge in volume, yet growing exponentially with time. It is a data with so large size and complexity that none of traditional data management tools can store it or process it efficiently. Big data is also a data but with huge size.



Fig 1 : Ten Big Data Applications in Real Life

Ref link : <https://images.app.goo.gl/QLxkM7gwFj963wus6>

Sources of Big Data

These data come from many sources like

- **Social networking sites:** Facebook, Google, LinkedIn all these sites generates huge amount of data on a day to day basis as they have billions of users worldwide.
- **E-commerce site:** Sites like Amazon, Flipkart, Alibaba generates huge amount of logs from which users buying trends can be traced.
- **Weather Station:** All the weather station and satellite gives very huge data which are stored and manipulated to forecast weather.

- **Telecom company:** Telecom giants like Airtel, Vodafone study the user trends and accordingly publish their plans and for this they store the data of its million users.
- **Share Market:** Stock exchange across the world generates huge amount of data through its daily transaction.

Big Data has played a pivotal role in the business environment today. We can understand this by looking at the aspects enlisted below,

- **Cost Savings:** Some tools of Big Data like Hadoop and Cloud-Based Analytics convey cost favorable circumstances to businesses when a lot of information is to be put away and these tools additionally help in distinguishing more proficient methods for working together.
- **Time Reduction:** The speedy nature of tools like Hadoop and in-memory analytics can undoubtedly recognize new sources of data which helps organizations in breaking down information instantly and identifying the most suitable decision.
- **Comprehend the economic situations:** Dissecting Big Data gives a clearer picture of the current economic scenario. For instance, by breaking down clients' buying practices, an organization can discover the items that are sold the most and deliver items as per this pattern. By this, it can stretch out beyond its rivals.

THE HISTORY OF BIG DATA

Although the concept of big data itself is relatively new, the origins of large data sets go back to the 1960s and '70s when the world of data was just getting started with the first data centers and the development of the relational database.

Around 2005, people began to realize just how much data users generated through Facebook, YouTube, and other online services. Hadoop (an open-source framework created specifically to store and analyze big data sets) was developed that same year. NoSQL also began to gain popularity during this time.

The development of open-source frameworks, such as Hadoop (and more recently, Spark) was essential for the growth of big data because they make big data easier to work with and cheaper to store. In the years since then, the volume of big data has skyrocketed. Users are still generating huge amounts of data—but it's not just humans who are doing it.

With the advent of the Internet of Things (IoT), more objects and devices are connected to the internet, gathering data on customer usage patterns and product performance. The emergence of machine learning has produced still more data.

While big data has come far, its usefulness is only just beginning. Cloud computing has expanded big data possibilities even further. The cloud offers truly elastic scalability, where developers can simply spin up ad hoc clusters to test a subset of data. And graph databases are becoming increasingly important as well, with their ability to display massive amounts of data in a way that makes analytics fast and comprehensive.

BIG DATA CHALLENGES

While big data holds a lot of promise, it is not without its challenges.

First, big data is...big. Although new technologies have been developed for data storage, data volumes are doubling in size about every two years. Organizations still struggle to keep pace with their data and find ways to effectively store it.

But it's not enough to just store the data. Data must be used to be valuable and that depends on curation. Clean data, or data that's relevant to the client and organized in a way that enables meaningful analysis, requires a lot of work. Data scientists spend 50 to 80 percent of their time curating and preparing data before it can actually be used.

Finally, big data technology is changing at a rapid pace. A few years ago, Apache Hadoop was the popular technology used to handle big data. Then Apache Spark was introduced in 2014. Today, a combination of the two frameworks appears to be the best approach. Keeping up with big data technology is an ongoing challenge.

Big data best practices

Here are some guidelines for building a successful big data foundation.

Align big data with specific business goals

More extensive data sets enable you to make new discoveries. To that end, it is important to base new investments in skills, organization, or infrastructure with a strong business-driven context to guarantee ongoing project investments and funding. To determine if you are on the right track, ask how big data supports and enables your top business and IT priorities. Examples include understanding how to filter web logs to understand ecommerce behavior, deriving sentiment from social media and customer support interactions, and understanding statistical correlation methods and their relevance for customer, product, manufacturing, and engineering data.

Ease skills shortage with standards and governance

One of the biggest obstacles to benefiting from your investment in big data is a skills shortage. You can mitigate this risk by ensuring that big data technologies, considerations, and decisions are added to your IT governance program. Standardizing your approach will allow you to manage costs and leverage resources. Organizations implementing big data solutions and strategies should assess their skill requirements early and often and should proactively identify any potential skill gaps. These can be addressed by training/crosstraining existing resources, hiring new resources, and leveraging consulting firms.

Optimize knowledge transfer with a center of excellence

Use a center of excellence approach to share knowledge, control oversight, and manage project communications. Whether big data is a new or expanding investment, the soft and hard costs can be shared across the enterprise. Leveraging this approach can help increase big data capabilities and overall information architecture maturity in a more structured and systematic way.

Top payoff is aligning unstructured with structured data

It is certainly valuable to analyze big data on its own. But you can bring even greater business insights by connecting and integrating low density big data with the structured data you are already using today.

Whether you are capturing customer, product, equipment, or environmental big data, the goal is to add more relevant data points to your core master and analytical summaries, leading to better conclusions. For example, there is a difference in distinguishing all customer sentiment from that of only your best customers. Which is why many see big data as an integral extension of their existing business intelligence capabilities, data warehousing platform, and information architecture.

Plan your discovery lab for performance

Discovering meaning in your data is not always straightforward. Sometimes we don't even know what we're looking for. That's expected. Management and IT needs to support this "lack of direction" or "lack of clear requirement." At the same time, it's important for analysts and data scientists to work closely with the business to understand key business knowledge gaps and requirements. To accommodate the interactive exploration of data and the experimentation of statistical algorithms, you need highperformance work areas. Be sure that sandbox environments have the support they need—and are properly governed.

Align with the cloud operating model

Big data processes and users require access to a broad array of resources for both iterative experimentation and running production jobs. A big data solution includes all data realms including transactions, master data, reference data, and summarized data. Analytical sandboxes should be created on demand. Resource management is critical to ensure control of the entire data flow including pre- and post-processing, integration, in database summarization, and analytical modeling.

A wellplanned private and public cloud provisioning and security strategy plays an integral role in supporting these changing requirements.

Use case

An e-commerce site XYZ (having 100 million users) wants to offer a gift voucher of 100\$ to its to

top 10 customers who have spent the most in the previous year. Moreover, they want to find the buying trend of these customers so that company can suggest more items related to them.

Issues

Huge amount of unstructured data which needs to be stored, processed and analyzed.

Solution

Storage: This huge amount of data, Hadoop uses HDFS (Hadoop Distributed File System) which uses commodity hardware to form clusters and store data in a distributed fashion. It works on Write once, read many times principle.

Processing: Map Reduce paradigm is applied to data distributed over network to find the required output.

Analyze: Pig, Hive can be used to analyze the data.

Cost: Hadoop is open source so the cost is no more an issue.

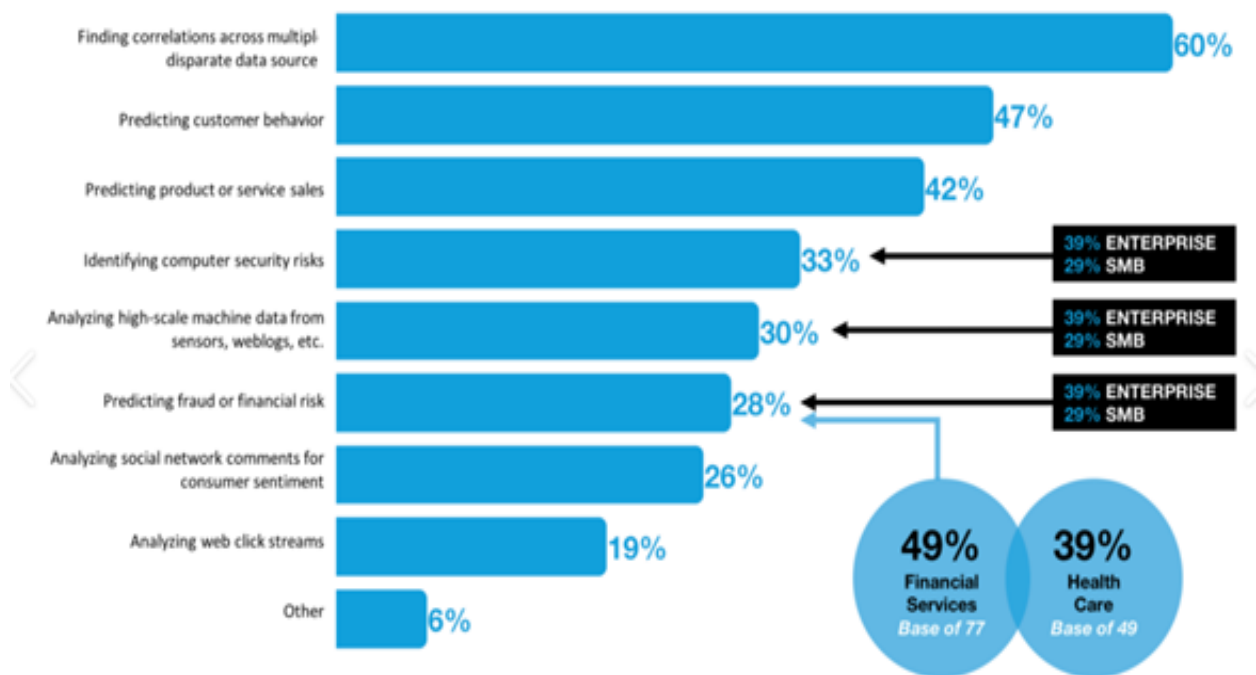


FIG 2 : Finding correlations across multiple disparate data source Predicting customer behavior.

Reflink:https://www.bing.com/images/search?view=detailV2&insightstoken=bcid_S1MdcS.JDlkDZkX7ZQ.0LmzoC1jg.....2g*ccid_Ux1xL8kO&form=SBIWEB&iss=SBIUPLOADGET&sbisrc=ImgPicker&idpbck=1&sbifsz=584+x+292+%c2%b7+21.54+kB+%c2%b7+png&sbifnm=Untitled.png&thw=584&thh=292&ptime=24&dlen=29404&expw=584&expw=292&selectedindex=0&id=1317468968&ccid=Ux1xL8kO&vt=3&sim=11&cal=0.05&cab=0.95&cat=0.05&car=0.95

Distributed Computing

A distributed computer system consists of multiple software components that are on multiple computers, but run as a single system. The computers that are in a distributed system can be physically close together and connected by a local network, or they can be geographically distant and connected by a wide area network. A distributed system can consist of any number of possible configurations, such as mainframes, personal computers, workstations, minicomputers, and so on. The goal of distributed computing is to make such a network work as a single computer.

Distributed systems offer many benefits over centralized systems, including the following:

Scalability

The system can easily be expanded by adding more machines as needed.

Redundancy

Several machines can provide the same services, so if one is unavailable, work does not stop. Additionally, because many smaller machines can be used, this redundancy does not need to be prohibitively expensive.

Distributed computing systems can run on hardware that is provided by many vendors, and can use a variety of standards-based software components. Such systems are independent of the underlying software. They can run on various operating systems, and can use various communications protocols. Some hardware might use UNIX or Linux as the operating system, while other hardware might use Windows operating systems. For intermachine communications, this hardware can use SNA or TCP/IP on Ethernet or Token Ring.

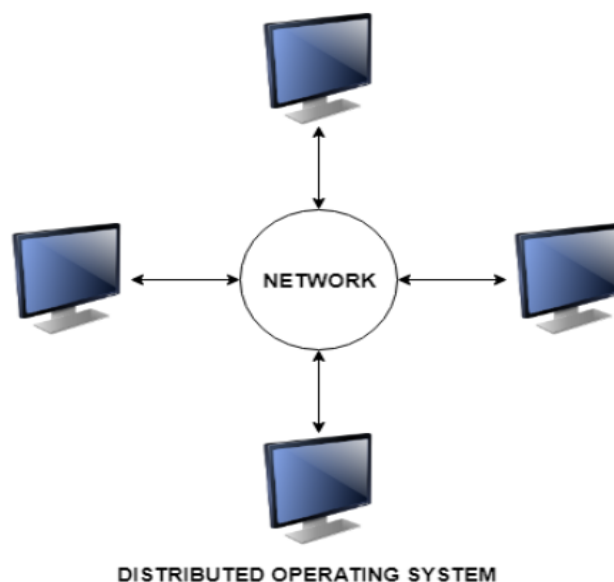


FIG: Distributed Operating System

Reflink:<https://www.tutorialspoint.com/Distributed-Systems>

In a distributed database management system, the database is not stored at a single location. Rather, it may be stored in multiple computers at the same place or geographically spread far away. Despite all this, the distributed database appears as a single database to the user.

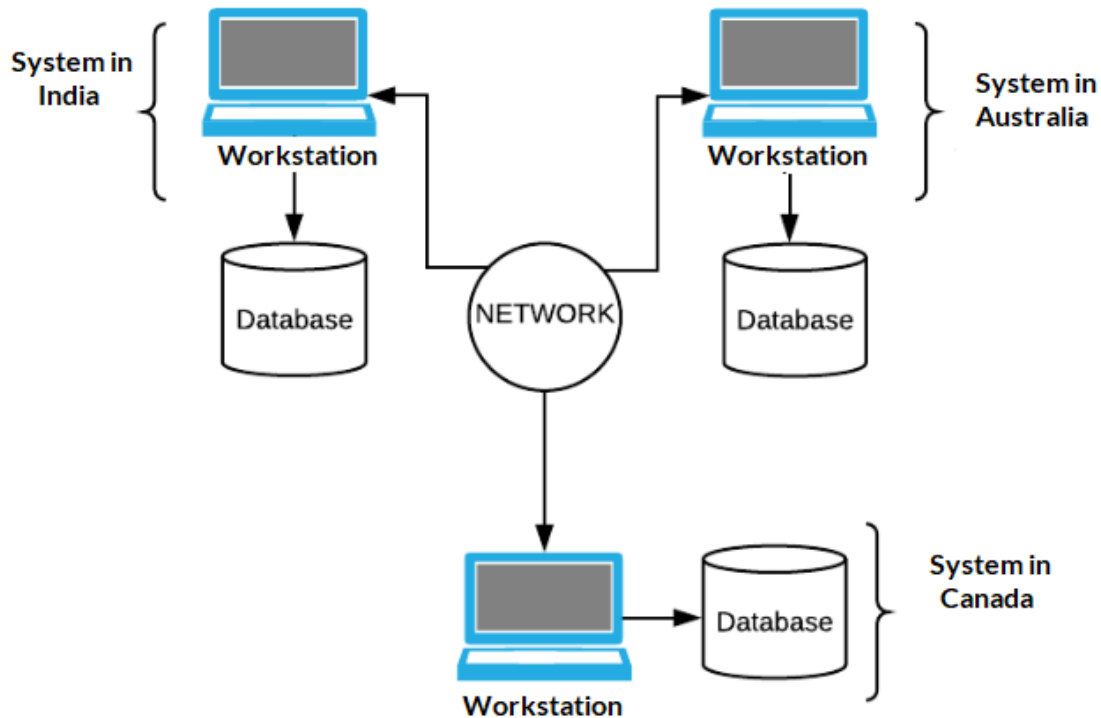


Fig:Network of different countries data bases

As seen in the figure, the components of the distributed database can be in multiple locations such as India, Canada, Australia, etc. However, this is transparent to the user i.e the database appears as a single entity.

Big Data with Distributed Computing Frame work

Distributed Computing has a great role in the success of Big Data. Big Data requires very low costing storage space and infrastructure, which is provided by cloud computing. Cloud Computing is a branch of Distributed Computing . In order to process a huge quantity of data at a very high speed, we required the power of cluster computing, which is also a branch of Distributed Computing. Thus, to enhance the processing speed of Big Data, there are two features: batch processing and stream processing . By default, Hadoop MapReduce is a batch processing system, but with the invent of social media inception of data in real-time, there is a need for a stream processing framework for Big Data. MapReduce processing frame work is designed to handle large data sets and split datasets into small batches. In contrast other frameworks process data in an uninterrupted stream, as it flows into the system.

MapReduce workings as its default processing engine but, for instance, Apache Spark and other framework are hooked into the Hadoop Ecosystem to handle real-time streaming data.

Comparison of different computing techniques considering different function

Grid Computing	Utility Computing	Cluster Computing	Cloud Computing
Loosely coupled	On-demand pricing	Tightly coupled systems	On-demand self-service
Diversity and dynamism	Uniform utility computing services	Single system image	Broad network access
Distributed job management and scheduling	Share the resources in the shared pool of machines	Centralized job management and scheduling system	Resources pooling and rapid elasticity
High-end computers (servers, clusters)	High-end computers (servers)	Commodity computers	Commodity computers, high-speed network and high-end servers and NAS

Batch processing frameworks split the large data jobs into small chunks and distribute these chunks on a large number of nodes according to the size of computer cluster to process Big Data. The execution time of a batch processing is determined with the number of active nodes in a cluster and the size of the job. The batch processing model is inappropriate to satisfy real-

time constraints due to having high latency to process Big Data. Stream processing is a model for handling real-

time stream synchronized with the data flow and returns the results in a low-

latency fashion. Stream processing also has some features of batch processing such as fault tolerance, high availability, and resource utilization. Real-

time stream processing systems give guarantee to be up and available all the time for real-

time data. Stream processing achieves incremental scalability automatically by distributing processing power as well as storage capacity across multiple computers without any human interaction. The following frameworks are hooked in Hadoop environment:

- MapReduce (Batch Processing frameworks)
- Spark (Stream and Batch Processing frameworks)
- Flink (Stream and Batch Processing frameworks)
- Storm (Stream processing frameworks)
- Samza (Stream Processing frameworks)

RELATED WORKS

Jun Ni, Ying Chen, Jie Sha, and Minghuan Zhang, Presented review on the demands and application potentials using big data technology with an emphasis on common challenges. After briefly addressing the Hadoop/MapReduce code components and modules, we use a simple clinic data to demonstrate how to map and reduce on small dataset with illustrated workflow. We give simple scenario of using other MapReduce calculation modules for counting and classification. This serves as a basic step into future utilization of big data to healthcare domain.

Yaping Chi , Yintan Yang, Ping Xu, Gefei Li, Shuhao Li, Presented review on addressing the limitations of Software Defined Networking (SDN) offers the means to dynamically configure the network parameters, dynamically provision network itself can be sliced in an on-demand manner. This research aims to characterize SDN with respect to the demands of big data analytics in Cluster, Grid, and Cloud Computing resources. The main motivation behind this research study is to design and develop an intelligent framework named as Big Data Analytics Management System (BDAMS) for collectively managing the compute, storage, and network resources in Cluster, Grid, and Cloud infrastructure for big data analytics.

J.Lozano, N.Aginako, M.Quartulli, I.Olaizola, E.Zulueta, P. Iriondo Presented review on contribution of describing a parallelized data processing approach for EO image analysis that is based on the MapReduce paradigm and implemented on the Apache Spark framework. Existing algorithms for e.g. thematic mapping need to be re-defined in order to exploit distributed execution capabilities to run on large coverage data.

Mohammed S. Al-

kahtani¹ Lutful Karim² and Jalal Almhana Presented review on CEDA for big data processing based on a framework that comprises data processing both at the data collection end and data processing server end. The proposed CEDA algorithm is application independent and scalable, i.e., using as many nodes as necessary. The proposed CEDA scheme supports both parallel and sequential implementation based on the amount of data to process. If the amount of data exceeds a certain threshold, it works in parallel mode. Otherwise, it works in sequential mode to reduce processing overhead. Moreover, the algorithm works on different types of nodes ranging from low powered RFID tags and sensors to high speed/powerful computers. The performance of the CEDA scheme was evaluated in terms of data size and data processing time. Simulation results show that the data size processed at the central server using CEDA is much smaller compared to that processed by existing approaches

Ivan E. Villalon-Turrubiates presented review paper on An intelligent post-processing paradigm based on the use of a dynamical filtering technique modified to enhance the reconstruction quality of remote sensing indexes using multitemporal images and distributed computing techniques is proposed. As a matter of particular study, a robust algorithm is reported for the analysis of the dynamic behavior of geophysical signatures extracted from remotely sensed scenes. The simulation results prove the efficiency of the proposed technique along with the computational implementation based on a big-data framework using distributed Processing.

Sandeep Kumar Hegde , Dr. Srinivasa K.G Presented review paper on Assigning the task to the several cluster of node in Map Reduce is an interesting problem, because efficient task assignment can significantly reduce the processing runtime, or improve hardware utilization. Many research work are proposed in Map Reduce on resource management to improve the performance of Hadoop with respect to the user of the application. A fair schedule strate

gy provides a delay Scheduling between jobs in pool to improve data locality. The heterogenous nodes are balanced using load balancing algorithm. In a heterogeneous computing environment the performance of Hadoop framework may be reduced due to lack of load management. In order to optimally allocate resources The taskbased algorithms and workflow based algorithms are used. For the data intensive application the Workflowbased approaches are implemented which will work even when the estimation about future tasks are not accurate but the problem with this technique is it is not suitable when load in the distributed environment is not uniform. The job scheduling was also done using Bayesian classification algorithm and the characteristic of this technique was it works based on the principle of Bayesian probability hence result is not optimal with this approach. To handle issue of local data the concept of wait scheduling is used and length of waitingtime is used as logic in order to schedule the data with time. Here the task are executed selectively which is not ordered strictly with respect to queue, so the problem of local data was improved. To improve the fairness of the task in Map Reduce cluster ,weight fair queuing scheduling algorithm is proposed . To allocate a weight to each and every queue and to schedule tasks of the sub- queue Weighted.

Alfredo Cuzzocrea and Ernesto Damiani presented review paper on theoretical data-driven privacy preserving big data framework in distributed environments can be designed, proved and extended, it is mandatory to set the different target data domains where the framework applies. Indeed, numerous specialized “vertical” realizations of the framework are possible, each one for each particular data setting. Among others, in this paper we consider the multidimensional data case because of, not only multidimensional data arise in a wide spectrum of relevant occurrences , but also they wellmarry with actual emerging big data analytics tools and systems where multidimensional analysis e.g., OLAP plays a major role.

Jun Ni, Ying Chen, Jie Sha, and Minghuan Zhang Presented review on the demands and application potentials using big data technology with an emphasis on common challenges . After briefly addressing the Hadoop/MapReduce code components and modules, we use a simple clinic data to demonstrate how to map and reduce on small dataset with illustrated workflow. We give simple scenario of using other MapReduce calculation modules for counting and classification. This serves as a basic step into future utilization of big data to healthcare domain

Tian-en Huang and Qinglai Guo Presented review on a framework for a distributed computing platform is designed. Then, distributed algorithms are developed, including a distributed massive sampling simulation method and a distributed feature selection method. Next, the software platform and hardware platform for the distributed computing platform are established. Finally, the platform is applied to the Guangdong Province Power System in China to evaluate its accuracy and efficiency. The simulation results show that the distributed computing platform can improve computing efficiency and perform better than a centralized platform.

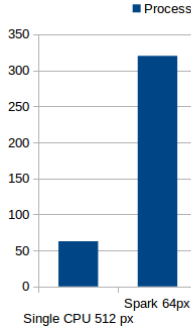
Lin Gu, Deze Zeng, Song Guo, Yong Xiang, Presented review on Big data stream processing (BDSP), has become a crucial requirement for many scientific and industrial applications in recent years. By offering a pool of computation, communication and storage resources,

public clouds, like Amazon's EC2, are undoubtedly the most efficient platforms to meet the ever-

growing needs of BDSP. Public cloud service providers usually operate a number of geographically distributed datacenters across the globe. Different datacenter pairs are with different inter-datacenter network costs charged by Internet Service Providers (ISPs). While, inter-datacenter traffic in BDSP constitutes a large portion of a cloud provider's traffic demand over the Internet and incurs substantial communication cost, which may even become the dominant operational expenditure factor.

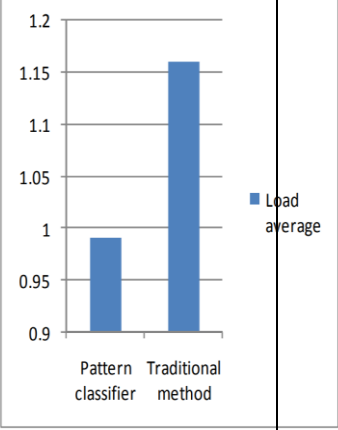
OVERVIEW OF ALL THE TECHNIQUES USED FOR THE DETECTION AND ANALYSIS OF BIG DATA AND DISTRIBUTED COMPUTING

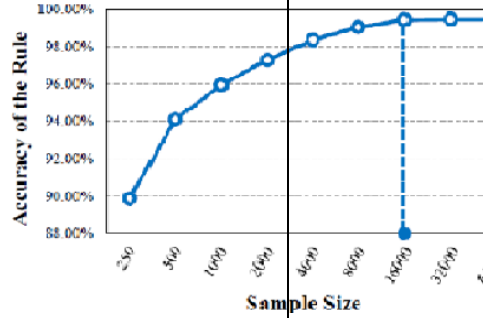
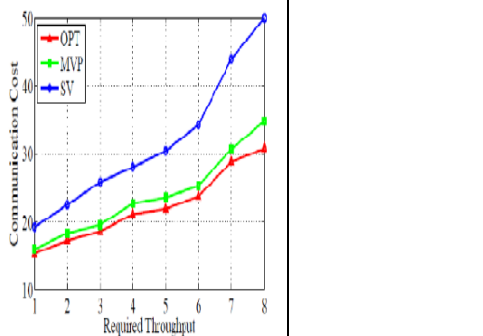
Sn o	Authors	Algorithm/Technique	Data sets	Accura cy	Result																										
1	Jun Ni, Ying Chen, Jie Sha, and Ming huan Zhang	Hadoop/MapReduce big data processing technique	Student Clinical Data	237.2	MapReduce mechanism <table border="1"> <thead> <tr> <th>StudentNo</th> <th>Clinic Fee</th> </tr> </thead> <tbody> <tr><td>13</td><td>12.80</td></tr> <tr><td>14</td><td>6.40</td></tr> <tr><td>15</td><td>6.40</td></tr> <tr><td>16</td><td>6.40</td></tr> <tr><td>17</td><td>12.80</td></tr> <tr><td>18</td><td>12.80</td></tr> <tr><td>19</td><td>5.60</td></tr> <tr><td>20</td><td>0.00</td></tr> <tr><td>21</td><td>6.40</td></tr> <tr><td>22</td><td>5.60</td></tr> <tr><td>23</td><td>6.40</td></tr> <tr><td>24</td><td>12.80</td></tr> </tbody> </table>	StudentNo	Clinic Fee	13	12.80	14	6.40	15	6.40	16	6.40	17	12.80	18	12.80	19	5.60	20	0.00	21	6.40	22	5.60	23	6.40	24	12.80
StudentNo	Clinic Fee																														
13	12.80																														
14	6.40																														
15	6.40																														
16	6.40																														
17	12.80																														
18	12.80																														
19	5.60																														
20	0.00																														
21	6.40																														
22	5.60																														
23	6.40																														
24	12.80																														
2	Yaping Chi , Yintan Yang,Ping Xu,Ge fei Li,Shuhao Li	MapReduce Processing Algorithm	traffic logs data	5GB Standalone Processing 363287 Ms Distributed Processing 20256 ms	COMPARISON OF STYLES DALONE PROCESSING AND DISTRIBUTED PROCESSING PERFORMANCE. HBase read performance test <table border="1"> <thead> <tr> <th>Data Way</th> <th>0.2GB</th> <th>0.5GB</th> <th>1GB</th> <th>3GB</th> </tr> </thead> <tbody> <tr> <td>Standalone processing</td> <td>26459ms</td> <td>61049ms</td> <td>150547ms</td> <td>239785ms</td> </tr> <tr> <td>Distributed processing</td> <td>35378ms</td> <td>54397ms</td> <td>139045ms</td> <td>181145ms</td> </tr> </tbody> </table>	Data Way	0.2GB	0.5GB	1GB	3GB	Standalone processing	26459ms	61049ms	150547ms	239785ms	Distributed processing	35378ms	54397ms	139045ms	181145ms											
Data Way	0.2GB	0.5GB	1GB	3GB																											
Standalone processing	26459ms	61049ms	150547ms	239785ms																											
Distributed processing	35378ms	54397ms	139045ms	181145ms																											

3	<p>J.Lozano, N.Aginako, M.Quartulli, I.Olaizola, E.Zulueta, P. Iriondo</p>	<p>Gaussian naive Bayesian supervised classification algorithm</p>	<p>original large scale dataset images</p>	<p>128 by 128 pixel sized tiles show</p>	 <table border="1"> <caption>Processing Time Comparison</caption> <thead> <tr> <th>Configuration</th> <th>Processing Time (Approx.)</th> </tr> </thead> <tbody> <tr> <td>Single CPU 512 px</td> <td>60</td> </tr> <tr> <td>Spark 64px</td> <td>320</td> </tr> </tbody> </table>	Configuration	Processing Time (Approx.)	Single CPU 512 px	60	Spark 64px	320
Configuration	Processing Time (Approx.)										
Single CPU 512 px	60										
Spark 64px	320										
4	<p>Mohammed S. Al-kahtani Lutful Karim2 and Jalal Almhana</p>	<p>Dyanamic Distributed Algorithm</p>	<p>Number of clusters at the data collecting side Max. 10 Number of collecting nodes in each cluster at data collecting side Max. 1024 Number of grouping nodes in each cluster at data collecting side Max. 20 Number of commodity computers Max. 60</p>	<p>results also showed that the amount of data to be processed by the central server and data processing time in the CEDDA approach</p>	<p>Simulation parameters and their values</p> <p>Parameter Value</p> <p>Number of clusters at the data collecting side Max. 10 Number of collecting nodes in each cluster at data collecting side</p>						

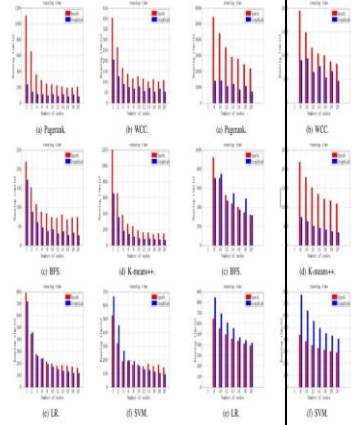
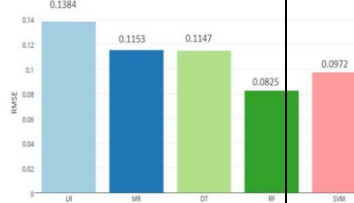
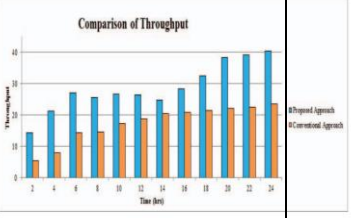
			<p>Size of data collected by a node at data collecting side 128 - 256 Bytes</p> <p>Duration of time data collected at collecting side 60 Minutes</p> <p>Number of groups at central server side 4</p> <p>Processing time of each data byte by commodity computer 0.25-0.5 nanosecond</p>		<p>Max 1024</p> <p>Number of grouping nodes in each cluster at data collecting side</p> <p>Max. 20</p> <p>Number of commodity computers</p> <p>Max. 60</p> <p>Size of data collected by a node at data collecting side 128 – 256 Bytes</p>
--	--	--	--	--	--

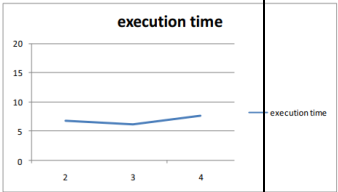
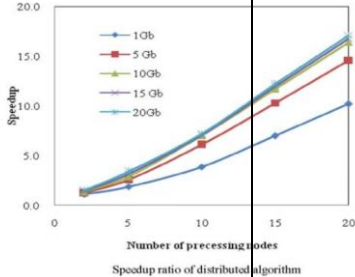
5	Ivan E. Villalon-Turrubiates IT in Industry, Vol. 9, No.3, 2021	Robust algorithm	distributed processing of large data sets across clusters of computers Published Online 30-Dec-2021	MDA framework in order to increase its accuracy is a matter	A RSS map is extracted from a set of multitemporal SPOT-5 images using the WPS method ie Black regions, Dark-
					<p>Duration of time data collected at collecting side</p> <p>60 Minutes</p> <p>Number of groups at central server side</p> <p>4</p> <p>Processing time of each data byte by commodity Computer</p> <p>0.25-0.5</p> <p>nanoseconds</p>

					gray regions, Light-gray regions						
6	SandeepKumar Hegde ,Dr. Srinivasa K.G	new pattern classifier algorithm	<p>Parameter For Running Job With Hadoop Cluster Parameter</p> <p>Description</p> <p>Job arrival distribution Exponential</p> <p>Job arrival rate 2 minutes</p> <p>Job Tracker Heart Beat Interval 1 Seconds</p> <p>Admission Interval 1 minutes</p> <p>Job Tracker map slots 60</p> <p>Job tracker Reduce slots 30</p>	Accuracy of computation	 <p>Comparison on Load Achieved</p> <table border="1"> <thead> <tr> <th>Method</th> <th>Load average</th> </tr> </thead> <tbody> <tr> <td>Pattern classifier</td> <td>~0.98</td> </tr> <tr> <td>Traditional method</td> <td>~1.15</td> </tr> </tbody> </table>	Method	Load average	Pattern classifier	~0.98	Traditional method	~1.15
Method	Load average										
Pattern classifier	~0.98										
Traditional method	~1.15										
7	Alfredo Cuzzocrea and Ernesto Damiani	privacy preserving big Multidimensional algorithm	privacy of big data sets sources		Emerging big multidimensional Data						
8	Jun Ni, Ying Chen, Jie Sha, and Minghuan Zhang	Traffic statistics analysis algorithm flow based on MapReduce		accuracy of the overall analytics process	DRIPROM, along with experimental evaluation and analysis; (ii) focusing on other types of emerging (big) data domains such as graph-like data and textual data						

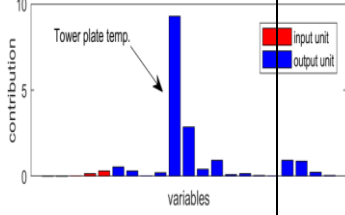
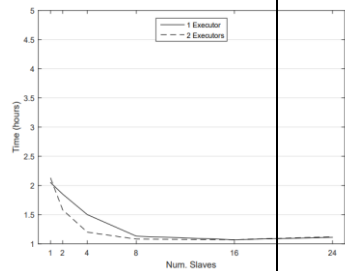
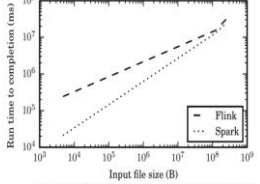
<p>9</p>	<p>Tian-en Huang and Qinglai Guo</p>	<p>distributed algorithms</p>	<p>Platform Model</p> <p>Centralized</p> <p>8-distributed</p> <p>Training Time (s)</p> <p>4,287</p> <p>754</p>	<p>Platform Model</p> <p>Centralized</p> <p>8-distributed</p> <p>Training Time (s)</p> <p>2200</p> <p>16000</p>	<p>The accuracy of the fine operational rule in scenarios with different sample sizes.</p>  <table border="1"> <caption>Accuracy of the Rule vs Sample Size</caption> <thead> <tr> <th>Sample Size</th> <th>Accuracy of the Rule (%)</th> </tr> </thead> <tbody> <tr><td>200</td><td>89.5</td></tr> <tr><td>800</td><td>94.0</td></tr> <tr><td>1600</td><td>96.0</td></tr> <tr><td>2400</td><td>97.5</td></tr> <tr><td>4000</td><td>98.5</td></tr> <tr><td>8000</td><td>99.0</td></tr> <tr><td>18000</td><td>99.5</td></tr> <tr><td>22000</td><td>99.0</td></tr> <tr><td>26000</td><td>99.0</td></tr> </tbody> </table>	Sample Size	Accuracy of the Rule (%)	200	89.5	800	94.0	1600	96.0	2400	97.5	4000	98.5	8000	99.0	18000	99.5	22000	99.0	26000	99.0																																
Sample Size	Accuracy of the Rule (%)																																																								
200	89.5																																																								
800	94.0																																																								
1600	96.0																																																								
2400	97.5																																																								
4000	98.5																																																								
8000	99.0																																																								
18000	99.5																																																								
22000	99.0																																																								
26000	99.0																																																								
<p>10</p>	<p>Lin Gu, Deze Zeng, Song Guo, Yong Xiang</p>	<p>Multiple Virtual Placement algorithm</p>	<table border="1"> <thead> <tr> <th>Local server</th> <th>Selectable linked server</th> <th>Distance from target</th> <th>Selected server</th> </tr> </thead> <tbody> <tr> <td>101</td> <td>(40, 134, 180)</td> <td>(182, 88, 42)</td> <td>180</td> </tr> <tr> <td>180</td> <td>(2, 101, 134, 200)</td> <td>(220, 121, 88, 22)</td> <td>200</td> </tr> <tr> <td>200</td> <td>(2, 180, 222)</td> <td>(220, 42, 0)</td> <td>222 (Target)</td> </tr> </tbody> </table>	Local server	Selectable linked server	Distance from target	Selected server	101	(40, 134, 180)	(182, 88, 42)	180	180	(2, 101, 134, 200)	(220, 121, 88, 22)	200	200	(2, 180, 222)	(220, 42, 0)	222 (Target)		 <table border="1"> <caption>Communication Cost vs Required Throughput</caption> <thead> <tr> <th>Required Throughput</th> <th>OPT</th> <th>MVP</th> <th>SV</th> </tr> </thead> <tbody> <tr><td>1</td><td>15</td><td>18</td><td>22</td></tr> <tr><td>2</td><td>18</td><td>22</td><td>28</td></tr> <tr><td>3</td><td>20</td><td>25</td><td>32</td></tr> <tr><td>4</td><td>22</td><td>28</td><td>38</td></tr> <tr><td>5</td><td>25</td><td>32</td><td>42</td></tr> <tr><td>6</td><td>28</td><td>35</td><td>48</td></tr> <tr><td>7</td><td>32</td><td>40</td><td>52</td></tr> <tr><td>8</td><td>35</td><td>45</td><td>55</td></tr> </tbody> </table>	Required Throughput	OPT	MVP	SV	1	15	18	22	2	18	22	28	3	20	25	32	4	22	28	38	5	25	32	42	6	28	35	48	7	32	40	52	8	35	45	55
Local server	Selectable linked server	Distance from target	Selected server																																																						
101	(40, 134, 180)	(182, 88, 42)	180																																																						
180	(2, 101, 134, 200)	(220, 121, 88, 22)	200																																																						
200	(2, 180, 222)	(220, 42, 0)	222 (Target)																																																						
Required Throughput	OPT	MVP	SV																																																						
1	15	18	22																																																						
2	18	22	28																																																						
3	20	25	32																																																						
4	22	28	38																																																						
5	25	32	42																																																						
6	28	35	48																																																						
7	32	40	52																																																						
8	35	45	55																																																						

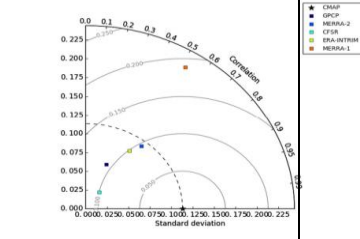
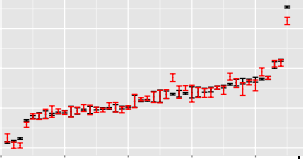
<p>11</p>	<p>A. Forestiero and G.Papuzzo</p>	<p>Distributed and Self organising algorithm.</p>	<table border="1"> <thead> <tr> <th>Local server</th> <th>Selectable linked server</th> <th>Distance from target</th> <th>Selected server</th> </tr> </thead> <tbody> <tr> <td>101</td> <td>(40, 134, 180)</td> <td>(182, 88, 42)</td> <td>180</td> </tr> <tr> <td>180</td> <td>(2, 101, 134, 200)</td> <td>(220, 121, 88, 22)</td> <td>200</td> </tr> <tr> <td>200</td> <td>(2, 180, 222)</td> <td>(220, 42, 0)</td> <td>222 (Target)</td> </tr> </tbody> </table>	Local server	Selectable linked server	Distance from target	Selected server	101	(40, 134, 180)	(182, 88, 42)	180	180	(2, 101, 134, 200)	(220, 121, 88, 22)	200	200	(2, 180, 222)	(220, 42, 0)	222 (Target)		<p>The average number of query hops performed by the algorithm to locate the target server. Simulation results, with solid lines, and mathematical results, with dot lines.</p>
Local server	Selectable linked server	Distance from target	Selected server																		
101	(40, 134, 180)	(182, 88, 42)	180																		
180	(2, 101, 134, 200)	(220, 121, 88, 22)	200																		
200	(2, 180, 222)	(220, 42, 0)	222 (Target)																		
<p>12</p>	<p>P. Li et al, Song Guo, Yoshiaki Miyazaki</p>	<p>Chance-constrained optimization technique.</p>	<p>Average performance of 30 MapReduce jobs</p> <p>SDA OPT taskOnly OPT exp OPT chance</p> <p>Inter-DC traffic 640 352.8 286.8 289.4</p> <p>Jobcompletion time 244.5 29.3172 31.4197 29.0284</p>	<p>55% inter-DC traffic</p>	<p>The average shuffling time under different number of reduce groups.</p>																
<p>13</p>	<p>YinanXu HuiLiu ZhihaoLong</p>	<p>A hybrid distributed computing framework</p>	<p>parallelized calculating of wind speed data</p>	<p>modified predictor on Spark</p>	<p>Improving percentages in MAE, MAPE and RMSE of IEWT reconstruction component on Spark can reach up to 38.778%, 39.438% and 37.227%, respectively</p>																

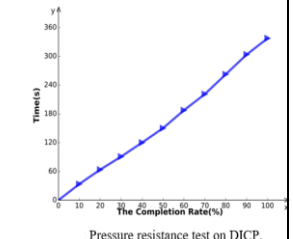
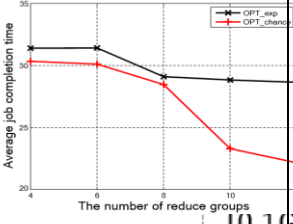
<p>14</p>	<p>J. Wei, K. Chen, Y. Zhou, Q. Zhou and J. He</p>	<p>Testing graph and generic analytic algorithms</p>			 <p>Fig. 1. Running time for algorithms with small dataset over local VM. Fig. 2. Running time for algorithms with large dataset over local VM.</p>
<p>15</p>	<p>Mohit Ved, Rizwanahmed B.</p>	<p>Random Forest algorithm</p>		<p>Multiple regression results in 62.6% accuracy with root mean squared error of 0.11531, while linear regression predicts with an accuracy of 40.2% with root mean squared error of 0.13839</p>	 <p>RMSE comparison of various Predictive Algorithms. LR: Linear Regression; MR: Multiple Regression; DT: Decision Tree; RF: Random Forest; SVM: Support Vector Machines</p>
<p>16</p>	<p>Kannan Govindarajan, Thamarai Selvi Somasundaram, David Boulanger, Vivekanandan Suresh Kumar, Kinshuk.</p>	<p>Software-defined Networking Technology, BigDataApplicationScheduler Algorithm</p>	<p>the submitted user requests</p>	<p>throughput of the submitted big data application requests.</p>	 <p>Comparison of Throughput</p>

<p>17</p>	<p>Zhi Yang ,Chunping Zhang ,Mu Hu ,Feng Lin</p>	<p>Objectification Parallel Computing (OPC).</p>	<p>Yarn based data sets</p>	<p>cluster changes from 2 to 4</p>	 <p>execution time</p> <p>Expansibility of system.</p>																									
<p>18</p>	<p>J. Yang, Y. Cao, B. Huang and Y. Zhao</p>	<p>Distributed Algorithm</p>	<p>Hiseq2500, Hiseq2500</p>	<p>greater than 5Gb and 10 processors.</p>	 <p>Speedup ratio of distributed algorithm</p>																									
<p>19</p>	<p>M. Cavallo, G. Di Modica, C. Polito and O. Tomarchio</p>	<p>SICP Mechanism</p>	<p>Geographically distant sites</p>	<p>test-bed was implemented to prove the viability of the approach test-bed was implemented to prove the viability of the approach test-bed was implemented to prove the viability of the approach</p>	<p>Table 7. Experiment results.</p> <table border="1" data-bbox="1742 839 2123 922"> <thead> <tr> <th>Data block location</th> <th>Global Reducer</th> <th>Real Execution Time [s]</th> <th>Predicted Execution Time [s]</th> <th>Error [%]</th> </tr> </thead> <tbody> <tr> <td>S_1, S_2, S_3</td> <td>S_2</td> <td>753</td> <td>698</td> <td>7%</td> </tr> <tr> <td>S_1, S_2, S_3</td> <td>S_2</td> <td>883</td> <td>812</td> <td>8%</td> </tr> <tr> <td>S_1, S_2, S_3</td> <td>S_2</td> <td>901</td> <td>818</td> <td>9%</td> </tr> <tr> <td>S_2, S_3, S_1</td> <td>S_2</td> <td>998</td> <td>911</td> <td>9%</td> </tr> </tbody> </table>	Data block location	Global Reducer	Real Execution Time [s]	Predicted Execution Time [s]	Error [%]	S_1, S_2, S_3	S_2	753	698	7%	S_1, S_2, S_3	S_2	883	812	8%	S_1, S_2, S_3	S_2	901	818	9%	S_2, S_3, S_1	S_2	998	911	9%
Data block location	Global Reducer	Real Execution Time [s]	Predicted Execution Time [s]	Error [%]																										
S_1, S_2, S_3	S_2	753	698	7%																										
S_1, S_2, S_3	S_2	883	812	8%																										
S_1, S_2, S_3	S_2	901	818	9%																										
S_2, S_3, S_1	S_2	998	911	9%																										

				Test bed proving the viability of the approach																																														
20	S. Dolev, P. Florissi, E. Gudes, S. Sharma and I. Singer	All Machine Learning algorithms	Geographical data sets	JetStream in SQL																																														
21	S. Bruce, Z. Li, H. Yang and S. Mukhopadhyay	A nonparametric two-sample inference algorithm																																																
22	P. Zhou, K. Wang, L. Guo, S. Gong and B. Zheng	Novel distributed federated online learning algorithm, T-PriDO Algorithm	predicting on <i>fixed-size</i> and small-scale datasets	<p>Average Accuracy.</p> <table border="1"> <thead> <tr> <th>Algorithm</th> <th>ϵ</th> <th>$n=1$</th> <th>$n=2$</th> <th>$n=3$</th> <th>$n=4$</th> <th>$n=5$</th> <th>$n=6$</th> <th>$n=7$</th> </tr> </thead> <tbody> <tr> <td>T-PriDO</td> <td>0.5</td> <td>41.92%</td> <td>44.89%</td> <td>47.72%</td> <td>54.37%</td> <td>55.91%</td> <td>58.43%</td> <td>60.26%</td> </tr> <tr> <td>T-PriDO</td> <td>2</td> <td>70.32%</td> <td>71.58%</td> <td>72.11%</td> <td>73.12%</td> <td>75.83%</td> <td>84.02%</td> <td>86.21%</td> </tr> <tr> <td>DT-PriDO</td> <td>0.5</td> <td>41.96%</td> <td>42.45%</td> <td>47.45%</td> <td>53.21%</td> <td>58.51%</td> <td>60.42%</td> <td>61.19%</td> </tr> <tr> <td>DT-PriDO</td> <td>2</td> <td>94.67%</td> <td>96.85%</td> <td>97.85%</td> <td>97.99%</td> <td>98.01%</td> <td>98.82%</td> <td>99.13%</td> </tr> </tbody> </table>	Algorithm	ϵ	$n=1$	$n=2$	$n=3$	$n=4$	$n=5$	$n=6$	$n=7$	T-PriDO	0.5	41.92%	44.89%	47.72%	54.37%	55.91%	58.43%	60.26%	T-PriDO	2	70.32%	71.58%	72.11%	73.12%	75.83%	84.02%	86.21%	DT-PriDO	0.5	41.96%	42.45%	47.45%	53.21%	58.51%	60.42%	61.19%	DT-PriDO	2	94.67%	96.85%	97.85%	97.99%	98.01%	98.82%	99.13%	<p>Average Regret-NT levels</p>
Algorithm	ϵ	$n=1$	$n=2$	$n=3$	$n=4$	$n=5$	$n=6$	$n=7$																																										
T-PriDO	0.5	41.92%	44.89%	47.72%	54.37%	55.91%	58.43%	60.26%																																										
T-PriDO	2	70.32%	71.58%	72.11%	73.12%	75.83%	84.02%	86.21%																																										
DT-PriDO	0.5	41.96%	42.45%	47.45%	53.21%	58.51%	60.42%	61.19%																																										
DT-PriDO	2	94.67%	96.85%	97.85%	97.99%	98.01%	98.82%	99.13%																																										
23	YinanXu HuiLiu Zhi haoLong	optimization algorithm	Resilient Distributed Datasets	performance of the modified predictor on Spark	<p>The forecasting results of Experiment 3 in Step 5.</p>																																													
24	J. Wei, K. Chen, Y. Zhou, Q. Zhou and J. He	generic analytic algorithms	VM's data sets	100% accuracy compared to over local computer clusters.																																														

25	Q. Jiang, S. Yan, H. Chen g and X. Yan	Distributed Modeling and Computing Framework for Nonlinear Process Monitoring.	Availability in industrial data	Optimal performance	 <p>Variable contribution plots of the distillation process fault 2.</p>
26	R. Talavera-Llames, R. Pérez-Chacón, A. Troncoso , F. Martínez-Álvarez	distributed algorithm	clustering techniques for small and medium datasets		 <p>Speed up depending on the number of slaves in a cluster.</p>
27	Milad Makkie, Xiang Li, Shannon Quinn, Binbin Lin, Jieping Ye, Geoffrey Mon, Tianming Liu	D-r1DL algorithm	tfMRI datasets	20% sampled data	 <p>Run time comparison of D-r1DL using Flink and Spark with varying input data sizes.</p>
28	Hu, J	Mathematical statistics and mining algorithms	All input variables	98.27%	

29	FeiHu^aChaoweiYang^a John L.Schnase^bDanielQ.Duffy^b MengchaoXu^aMichael K.Bowen^bTsengdarLee^cWeiwei Song^a	in-memory, distributed computing framework	Segments of data sets	0.98%	 <p>The Taylor diagram result for the monthly precipitation fr (-45°, -90°)~(45°, 90°)</p>															
30	Shiow-LuanWangYung-TsungHou	Approximation algorithm	All Input variables	0.1 set of performance	<p>Results of work line reliability test.</p> <table border="1" data-bbox="1758 582 2235 678"> <thead> <tr> <th>Group</th> <th>A</th> <th>B</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.95</td> <td>0.95</td> </tr> <tr> <td>2</td> <td>0.95</td> <td>0.97</td> </tr> <tr> <td>3</td> <td>0.94</td> <td>0.95</td> </tr> <tr> <td>System reliability ($R = 1 - (1-R1) \times (1-R2) \times (1-R3)$)</td> <td>0.98</td> <td></td> </tr> </tbody> </table>	Group	A	B	1	0.95	0.95	2	0.95	0.97	3	0.94	0.95	System reliability ($R = 1 - (1-R1) \times (1-R2) \times (1-R3)$)	0.98	
Group	A	B																		
1	0.95	0.95																		
2	0.95	0.97																		
3	0.94	0.95																		
System reliability ($R = 1 - (1-R1) \times (1-R2) \times (1-R3)$)	0.98																			
31	Scott Bruce, , Zeda Li, Hsiang-Chieh Yang, SubhadeepMukhopadhyay	nonparametric two sample inference algorithm	statistical modeling tools of large datasets	95% confidence intervals																
32	Rohyoung Myung# , Heonchang Yu# , Daewon Lee	Machine learning algorithm	data analytics programming data sets	1.5 to 3.3 times improvement of execution time																
33	Ivan E. Villalon-Turrubiates	Big-Data Technique	data sets across clusters of computers																	

<p>34</p>	<p>Le Dong, Zhiyu Lin, Yan Liang, Ling He, Ning Zhang, Qi Chen, Xiaochun Cao and Ebroul Izquierdo</p>	<p>SICP & DICP Mechanism</p>	<p>ImageNet Dataset</p>		 <p>Pressure resistance test on DICP.</p>
<p>35</p>	<p>Peng Li, Song Guo, Toshiki Miyazaki, Xiaofei Liao,</p>	<p>linearization and relaxation algorithm</p>	<p>shuffling time of 30 MapReduce job instances</p>		 <p>The average job completion time</p> <p>The number of reduce groups</p> <p>OPT_exp</p> <p>OPT_chance</p> <p>The average</p> <p>The average</p> <p>The average</p> <p>[0.1GB, 0.6GB]</p> <p>B)</p>

					[0.1GB, 0.8GB] B [0.1GB, 1GB]
--	--	--	--	--	-------------------------------------

Conclusion

There is an increasing trend of analyzing big data distributed over several data centers located in different countries and regions. Big data has received significant attention from researchers, business industries, education, and scientific communities. It consists of both unstructured and structured data that should be properly extracted, processed, and analyzed in order to obtain meaningful information. It requires large amount of high performance compute cycles, storage, and network bandwidth. The proposed research work first aims to aggregate the high performance computing resources from cluster, grid, and cloud as a collective infrastructure using distributed algorithms. In future, we will continue to study the system implementation by integrating proposed algorithms into popular data processing platform

References

- [1] Roger Magoulas and Lorica Ben, "Big Data: Technologies and Techniques for Large-Scale Data", Release 2.11, O'Reilly Media Inc., 2011.
- [2] P. Zikopoulos and C. Eaton, "Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data," New York, NY, USA: McGraw-Hill, 2011.
- [3] H. Barwick. The "four Vs" of big data. Implementing Information Infrastructure Symposium. [Online]. Available: http://www.computerworld.com.au/article/396198/iiis_four_vs_big_data
- [4] Hong Bin, Peng Fuyang, Deng Bo, Wang Dongxia. A survey on state monitoring of computational resources in cloud[J]. Computer Applications and Software, 2016(6).
- [5] Q. Wang, H. Wang, C. Zhang, W. Wang, Z. Chen and F. Xu, "A Parallel Implementation of Idea Graph to Extract Rare Chances from Big Data", *IEEE Intl Conference on Data Mining Workshop*, Shenzhen, 2014, pp. 503-510.
- [6] J. P. Verma, B. Patel, and A. Patel, "Big data analysis: recommendation system with Hadoop framework", *IEEE Proc. of CICT*, pp. 92-97, 2015.
- [7] Crespo, A...Garcia-Molina, H: Routing indices for peer-to-peer systems. In: Proc. of the 22nd International conference on Distributed Computing systems ICDCS '02. pp. 23-33 (2002)
- [8] S. K. Hegde and Srinivasa K.G, "A novel pattern classifier approach towards the performance optimization of Big Data analysis in distributed environment," 2017 Third International Conference on Advances in Electrical, Electronics, Information, Communication and Bio-Informatics (AEEICB), 2017, pp. 100-104, doi: 10.1109/AEEICB.2017.7972391.
- [9] A. Cuzzocrea and E. Damiani, "Pedigree-ing Your Big Data: Data-Driven Big Data Privacy in Distributed Environments," 2018 18th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2018, pp. 675-681, doi: 10.1109/CCGRID.2018.00100
- [10] Y. Chi, Y. Yang, P. Xu, G. Li and S. Li, "Design and implementation of monitoring data storage and processing scheme based on distributed computing," 2018 IEEE 3rd International Conference on Big Data Analysis (ICBDA), 2018, pp. 206-211, doi: 10.1109/ICBDA.2018.8367678.
- [11] J. Lozano, N. Aginako, M. Quartulli, I. Olaizola, E. Zulueta and P. Iriondo, "Large scale thematic mapping by supervised machine learning on 'big data' distributed cluster computing frameworks," 2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), 2015, pp. 1504-1507, doi: 10.1109/IGARSS.2015.7326065.

- [12] Z. Yang, C. Zhang, M. Hu and F. Lin, "An Electric Power Big Data Deployment Solution for Distributed Memory Computing," 2015 Seventh International Symposium on Parallel Architectures, Algorithms and Programming (PAAP), 2015, pp. 155-161, doi: 10.1109/PAAP.2015.38.
- [13] Z. Yang, C. Zhang, M. Hu and F. Lin, "OPC: A Distributed Computing and Memory Computing-Based Effective Solution of Big Data," 2015 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity), 2015, pp. 50-53, doi: 10.1109/SmartCity.2015.46.
- [14] S. Jain, A. Kumar, S. Mandal, and etc., "B4: experience with globally-deployed software defined wan," in *Proc. of ACM SIGCOMM*, Hong Kong, China, August 2013, pp. 3-14.
- [15] J. Yang, Y. Cao, B. Huang and Y. Zhao, "A Distributed Algorithm for Quality Assessment of Biological Sequencing Based on MapReduce," 2019 IEEE 5th International Conference on Computer and Communications (ICCC), 2019, pp. 188-192, doi: 10.1109/ICCC47050.2019.9064159.
- [16] N. Bahareva, Y. Ushakov, M. Ushakova, D. Parfenov, L. Legashev and I. Bolodurina, "Researching a Distributed Computing Automation Platform for Big Data Processing," 2020 International Conference Engineering and Telecommunication (En&T), 2020, pp. 1-5, doi: 10.1109/EnT50437.2020.9431254.
- [17] S. K. Hegde and Srinivasa K.G, "A novel pattern classifier approach towards the performance optimization of Big Data analysis in distributed environment," 2017 Third International Conference on Advances in Electrical, Electronics, Information, Communication and Bio-Informatics (AEEICB), 2017, pp. 100-104, doi: 10.1109/AEEICB.2017.7972391.
- [18] Y. Chi, Y. Yang, P. Xu, G. Li and S. Li, "Design and implementation of monitoring data storage and processing scheme based on distributed computing," 2018 IEEE 3rd International Conference on Big Data Analysis (ICBDA), 2018, pp. 206-1, doi:10.1109/ICBDA.2018.8367678.
- [19] L. Dong et al., "A Hierarchical Distributed Processing Framework for Big Image Data," in *IEEE Transactions on Big Data*, vol. 2, no. 4, pp. 297309, 1 Dec. 2016, doi: 10.1109/TBDDATA.2016.2613992
- [20] B. R. Chang, H. Tsai, Y. Wang and C. Huang, "Resilient distributed computing platforms for big data analysis using Spark and Hadoop," 2016 International Conference on Applied System Innovation (ICASI), 2016, pp. 1-4, doi: 10.1109/ICASI.2016.7539859.
- [21] M. Cavallo, G. Di Modica, C. Polito and O. Tomarchio, "H2F: A Hierarchical Hadoop Framework for Big Data Processing in GeoDistributed Environments," 2016 IEEE/ACM 3rd International Conference on Big Data Computing Applications and Technologies (BDCAT), 2016, pp. 27-35.
- [22] P. Zhou, K. Wang, L. Guo, S. Gong and B. Zheng, "A PrivacyPreserving Distributed Contextual Federated Online Learning Framework with Big Data Support in Social Recommender Systems," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 3, pp. 824-838, 1 March 2021, doi: 10.1109/TKDE.2019.2936565.
- [23] A. Forestiero and G. Papuzzo, "Distributed Algorithm for Big Data Analytics in Healthcare," 2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI), 2018, pp. 776-779, doi: 10.1109/WI.2018.00015
- [24] J. Lozano, N. Aginako, M. Quartulli, I. Olaizola, E. Zulueta and P. Iriondo, "Large scale thematic mapping by supervised machine learning on 'big data' distributed cluster computing frameworks," 2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), 2015, pp. 1504-1507, doi: 10.1109/IGARSS.2015.7326065

- [25] R. S. Gargees and G. J. Scott, "MultiStage Distributed Computing for Big Data: Evaluating Connective Topologies," 2020 10th Annual Computing and Communication Workshop and Conference (CCWC), 2020, pp. 0626-0633, doi: 10.1109/CCWC47524.2020.9031227
- [26] P. Li et al., "TrafficAware GeoDistributed Big Data Analytics with Predictable Job Completion Time," in IEEE Transactions on Parallel and Distributed Systems, vol. 28, no. 6, pp. 1785-1796, 1 June 2017, doi: 10.1109/TPDS.2016.2626285
- [27] Z. Yang, C. Zhang, M. Hu and F. Lin, "An Electric Power Big Data Deployment Solution for Distributed Memory Computing," 2015 Seventh International Symposium on Parallel Architectures, Algorithms and Programming (PAAP), 2015, pp. 155-161, doi: 10.1109/PAAP.2015.38
- [28] S. Dolev, P. Florissi, E. Gudes, S. Sharma and I. Singer, "A Survey on Geographically Distributed BigData Processing Using MapReduce," in IEEE Transactions on Big Data, vol. 5, no. 1, pp. 60-80, 1 March 2019, doi: 10.1109/TBDDATA.2017.2723473
- [29] S. Bruce, Z. Li, H. Yang and S. Mukhopadhyay, "Nonparametric Distributed Learning Architecture for Big Data: Algorithm and Applications," in IEEE Transactions on Big Data, vol. 5, no. 2, pp. 166-179, 1 June 2019, doi: 10.1109/TBDDATA.2018.2810187
- [30] M. Makkie et al., "A Distributed Computing Platform for fMRI Big Data Analytics," in IEEE Transactions on Big Data, vol. 5, no. 2, pp. 109119, 1 June 2019, doi: 10.1109/TBDDATA.2018.2811508.
- [31] L. Gu, D. Zeng, S. Guo, Y. Xiang and J. Hu, "A General Communication Cost Optimization Framework for Big Data Stream Processing in GeoDistributed Data Centers," in IEEE Transactions on Computers, vol. 65, no. 1, pp. 1929, 1 Jan. 2016, doi: 10.1109/TC.2015.2417566.
- [32] M. Ved and R. B., "Big Data Analytics in Telecommunication using State-of-the-art Big Data Framework in a Distributed Computing Environment: A Case Study," 2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC), 2019, pp. 411-416, doi: 10.1109/COMPSAC.2019.00066
- [33] J. Ni, Y. Chen, J. Sha and M. Zhang, "HadoopBased Distributed Computing Algorithms for Healthcare and Clinic Data Processing," 2015 Eighth International Conference on Internet Computing for Science and Engineering (ICICSE), 2015, pp. 188193, doi: 10.1109/ICICSE.2015.41
- [34] M. S. Alkahtani, L. Karim and J. Almhana, "Computationally Efficient, Dynamic distributed Algorithm of sensorbased Big Data," 2017 13th International Wireless Communications and Mobile Computing Conference (IWCMC), 2017, pp. 759-763, doi: 10.1109/IWCMC.2017.7986380
- [35] K. Govindarajan, T. S. Somasundaram, D. Boulanger, V. S. Kumar and Kinshuk, "A framework for scheduling and managing big data applications in a distributed infrastructure," 2015 Seventh International Conference on Advanced Computing (ICoAC), 2015, pp. 1-6, doi: 10.1109/ICoAC.2015.7562784.
- [36] A. Cuzzocrea and E. Damiani, "Pedigreeing Your Big Data: DataDriven Big Data Privacy in Distributed Environments," 2018 18th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2018, pp. 675681, doi: 10.1109/CCGRID.2018.00100
- Z. Yang, C. Zhang, M. Hu and F. Lin, "OPC: A Distributed Computing and Memory ComputingBased Effective Solution of Big Data," 2015 IEEE International Conference on Smart City/SocialCom/SustainCom (SmartCity), 2015, pp. 50-53, doi: 10.1109/SmartCity.2015.46
- J. Yang, Y. Cao, B. Huang and Y. Zhao, "A Distributed Algorithm for Quality Assessment of Biological Sequencing Based on MapReduce," 2019 IEEE 5th International Conference on Computer and Communications (ICCC), 2019, pp. 188-192, doi: 10.1109/ICCC47050.2019.9064159