

An Analytical Approach to Analyze Big Web Data

Sandeep¹, Sachin Kumar Chauhan², Reema³, Shabnam Sangwan⁴

^{1,2}Mtech scholar, Department of CSE, Sat Kabir Institute of Technology & Management, Bahadurgarh, Haryana, India
sk1034167@gmail.com

sonuchauhan63@gmail.com

^{3,4}A.P., Department of CSE, Sat Kabir Institute of Technology & Management, Bahadurgarh, Haryana, India
Shabnam022@email.com

arorareema@live.com

Abstract:- The basic of this dissertation is all about the rise of "big data" and the use of analytics to storehouse the data. Big data is basically used to analyze a large amount of data. The Big data store a large amount of data and provide helpful information in an efficient manner leads a system to serious computational challenges, like to analyze, mixture, & store, where information are remotely collected. In the recent times, data warehousing, data repository, and data mining are mostly used for Bid data. Big data warehouse known as terabytes which mean data collect in the warehouse was terabytes in storage but in recent time it is petabytes, and the data built at a high speed. The progress in companies and organization is mainly because it stores and analyzes the data at greater levels and in greater details, as well as metadata, Web data, and code generated data, to build a better relationship between customer and market behavior. Big data provide helpful information to achieve a goal. A lot of companies in recent time use big data to improve their quality. The growth of big data increases at a higher speed in recent time, the even trend going in a coming year. Keywords: big data, data warehouse, a data repository.

1. INTRODUCTION

This paper is in the space of data warehouse in the large industry to keep pace with the desire to collect and analyze large data or volumes of informational data and structured data. A data warehouse is a collection of a large amount of information as well as supporting system and this information is used to build project or software. The vendor of RDBMS has provided a different platform and each platform has a different specialty which provides a higher level of price replication and it also provides higher performance as compared to general purpose Relational DBMS. These platforms provide a lot of information is available in a variety of shapes and size, to the database. there is a large number of the survey is done by people which give real-time description .therefore, time to time updating is done by the company. We know that new technologies are coming to handle complex data. These technologies provide a large amount of complex data which include Metadata, web data and server data. The Metadata mean data about data which mainly provide a large amount of data. Metadata basically provide descriptive information about data. The Web data gives data about social media data like what Sapp data, messenger data, Facebook data etc. A large amount of data is uploaded to the Facebook server, the data in the Facebook server is very large and these types of data are social media content. New technologies provide machine generated data or code generated data and these technologies have also advance feature like sensors and data like GPS. The kind of data is known as big data. Big data mean very large data. These kinds of data have in large volume and these kinds of data are used to build or analyze

new technologies. Now the question arises what is new technologies and how we use these technologies world. Let discuss in brief. A large amount of company, which provide open framework it means any person can use these open framework and work on them. And these Data is change structured into unstructured data. These data was including in batch jobs that run on the server machine. The term BIG DATA is used to address the data set which mainly provides large and complex data. The traditional application is not enough to handle large and complex data. These mainly used for the purpose of analyzing, queues, capture, for the purpose sharing data, solve complex queries, search data for their application and protect information from the unauthorized user. This analytical method can be used to describe the web data. The basic advantage of Big Data is accuracy which leads your data more confidential in decision making. If decisions are better they can provide better performance, more efficient, effective and increase productivity and reduce risk with the low cost so it can be cost effective. Big data analytics is basically a process which is used to check or examine a bit of data to uncover a hidden pattern. It well known about the customer preference and provide a better result for their customer. Big data provide large data in small format. Data store in big data is in the range of gigabytes to terabytes and now to the petabytes. There is multiple data algorithm which provides information at each level in a specific manner and also in detail manner.

1.1 HADOOP: BIG DATA OVERVIEW

Hadoop is a key framework for the analysis of big data. It is an open-source framework that allows to store and process large volume data, structured and unstructured complex

data. This framework is designed to analyze very large data and also can be used to scale up from single servers to thousands of servers, with each server has the capacity to perform local computation and has its local storage. With the invention of new technologies, new devices and new and advance communication means like Facebook, Twitter, and WhatsApp, the data is growing in volume at a very high speed. Data produced by only social networking sites is the amount of data that if we bind up the data into the form of disks it might fill an entire football ground. This shows that social networking sites are producing the very large amount of data. The data produced by these sites was in gigabyte at the start of 21st century. But it grew at a very large speed. Then in terabytes and now in petabytes. In coming years this data is going to be very large due to the advent of new technologies and new communication devices. To process this large data we were requiring a framework. Then Hadoop came into picture which is open source. This framework is used to process and analyze the big data. With the help of Hadoop, we can run applications with a large number of nodes with large number data which may be in terabytes. This allows the very fast transfer of data between nodes. It is also a very reliable framework which allows continuing operating if even any node fails. It also minimizes the risk and chance of system failure even if a good number of nodes get failed.

1.2 MAPREDUCE: HEART OF BIG DATA

MapReduce is the heart of Big Data. MapReduce is also a framework which is used to write applications to process and analyze large amounts of data in parallel on big servers in a reliable manner. It is programming method or programming technique which allows for large scalability across a large number of servers, may be hundreds or even thousands. The MapReduce concept is very simple to understand. For those who are familiar with large scale-out data processing solutions, this concept is simpler. For new commerce, it is somewhat difficult to understand because it is not like the traditional system. MapReduce actually refers to two different tasks which are performed by Hadoop. The first task is the Map task, which takes input as a big data in the form of data sets and converts them into another set of data, in which elements are broken down into further lower elements called tuples. Tuples are the key/value pairs. The Second task is the Reduce task. Reduce task takes the output from a Map task as input and combines those datasets tuples into further smaller tuples. As the sequence of the name MapReduce implies, the reduce job is always performed after the map job.

The main advantage of MapReduce is its accuracy and also it is easier to process data over a large number of computing nodes. In the MapReduce model, the data processing

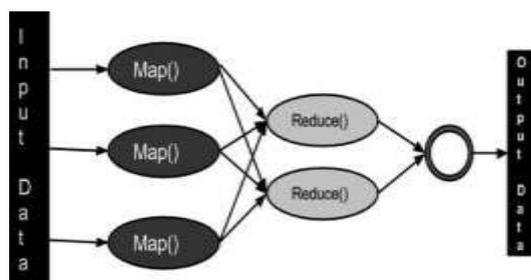
attributes are called mappers and reducers. MapReduce decomposes the data processing applications into Mappers and Reducers. This Decomposition of the application into mappers and reducers is sometimes non-trivial that it may have some variables or terms that are not equal to zero or that are not equal to identity. But, once we had written an application using MapReduce model that is we had decomposed the application into Mapper and Reducer then we can process or scale the application to run over a large number, may be hundreds, thousands or might be tens of thousands of machines in a cluster by changing only a few configurations. This is the main feature of MapReduce model which enables many programmers to use the MapReduce model. Generally, MapReduce model is based on sending the computer to the place where the actual data resides.

The execution of MapReduce program consists of three stages which are: Map stage, Shuffle stage and Reduce stage.

Map stage: The Map stage or Mapper's task is to process the big data which is taken as input to the mapper. This input data is in the form of datasets or file or directory which is stored in the Hadoop file system (HDFS). The input datasets are passed to the Mapper function line by line. The Mapper processes the data and decompose them into smaller elements or creates several small chunks of data.

Reduce stage: The Reduce stage or Reducer's task is to process the data which is output from the Mapper. After processing the data which comes from Mapper, Reducer produces a new set of output, which is then stored in the HDFS.

In this process of MapReduce, Hadoop sends the Map and Reduce tasks to the respective servers in the cluster. The Hadoop framework manages all the processing and details of data-passing such as issuing tasks, verifying task completion, and transfer the data around the cluster between the nodes. Most of the computation takes place on nodes with data on local disks due to which the traffic on the network gets reduced. After these tasks are completed, the framework collects the data and reduces it to form the required appropriate result, and then sends it back to the Hadoop server.



2. BIG DATA SIGNIFICANCE

Big Data is growing at a regular interval at a very high speed. It is becoming so large and complex that it is difficult to process, store, manage, share and analyze within current computational powers or with the traditional processing applications. Big data is a very important field. It is becoming the growing field. Almost all big organizations and social media companies are lying on Big Data. There are four key characteristics of big data which make them so important and significant. These characteristics are Volume, Velocity, and Variety. These are called 3vs of Big Data.

Volume: Volume is the most important characteristic of Big Data. As we all know, Web Data is growing in volume. With the help of Big Data frameworks available in the market like Hadoop and MapReduce, we can manage, store process and analyses such large volume of data easily.

Velocity: By the advent of Big Data it is data streaming at a high speed is now possible. This type of framework is real-time. For example, we can easily stream video on YouTube. Also, the larger volume of Twitter data ensured high velocity of even at 140 characters per tweet.

Variety: This characteristic refers to the fact that Big Data is capable of processing a variety of data which includes Structured, Unstructured, Semi-Structured, Text, Images and other media and so on. With the help of Big Data technology, we are now capable of analyzing and processing this variety of data.

Finally, social media sites like Facebook, Twitter and LinkedIn would not exist without the big data.

videos on the YouTube. By using, MapReduce it is possible to analyze the web data. The Analysis is described below.

DATA SET DESCRIPTION

The dataset or data table for YouTube consist of the following information.

- 1: A Unique Video id.
- 2: Details about Uploader of the video.
- 3: Time gap between the date of uploading of the video and YouTube establishment date.
- 4: Category to which video belongs to.
- 5: Length of time of the video.
- 6: Number of views till now.
- 7: Rating.
- 8: Number of ratings.
- 9: Number of comments.
- 10: Ids of related videos.

Problem Statement 1

Our motive is to find out the top 5 categories having a maximum number of videos uploaded.

SOURCE CODE

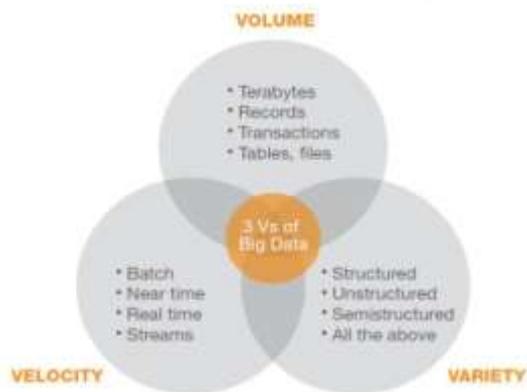
Now by using Mapper code, we will get the video category as key and final int value as values which will be passed to the next stage which is shuffle and sort phase and then this value is sent to the reducer phase where the aggregation of the values is performed.

Mapper Code

```
public class TopFiveCategories {
    public static class Map_1 extends
    Mapper<LongWritable, Text, Text,
    IntWritable>
    {
        private Text cat = new Text();
        private final static IntWritable one_1 = new IntWritable(1);
        public void map1(LongWritable key, Text value, Context
        context )
        throws IOException, InterruptedException
        {
            String one_line = value.toString();
            String str[]=one_line.split("\t");
            if(str.length > 5){
                cat.set(str[3]);
            }
            context.write(cat, First);
        }
    }
}
```

Explanation of the above Mapper code: In line 1 we are taking a class by name TopFiveCategories. In line 2 we are extending the Mapper default class with arguments KeyIN, ValueIN, KeyOut, ValueOut

In line 3 we are declaring a private Text variable 'cat' that is a category which will store the category of videos on YouTube.



3. PROPOSED WORK

This dissertation is about analyzing the data of YouTube. This analysis is performed using the Hadoop and MapReduce frameworks. The YouTube data is open source and publically available to the users. This YouTube data set is described below under the Data Set Description section. Using this dataset we will perform some Analysis and will find some results of this analysis like what are the top rated videos on YouTube also who uploaded the most number of

In line 4 we are declaring a variable which is private final static IntWritable variable one_1 which will be a constant value for every value.

In line 5 we are overriding the map_1 method.

In line 7 we are storing the line in a string variable one_line.

In line 8 the line is split using comma “,” delimiter and value are stored in an array of string so that all the columns in a row are stored in the string array.

In line 9 if the block is used to check whether the string array of length is greater than 6 which means it will enter into the if block and execute the code to eliminate the Exceptions.

In line 10 we are storing the cat that is a category.

In line 12 we are writing the key and. Which will be the output of the Mapper’s Map_1 method.

Reduce Code

```
public static class ReduceCat extends Reducer<Text, IntWritable,Text,IntWritable>
{
public void reducecat(Text key, Iterable<IntWritable> values,Context context throws IOException, InterruptedException
{
int sumcat = 0;
for (IntWritable val : values)
{
sumcat += val.get();
}
context.write(key, new IntWritable(sumcat));
}
}
```

While coming to the Reducer code:

line 1 a class Reduced cat which is extended the default Reducer class with arguments KeyIn, ValueIn, same as the outputs of the mapper class and, ValueOut that is used to store final outputs of our MapReduce program.

In line 2 Reduce method is overridden which is run each time for every key.

In line 3 An integer variable sum cart is declared which is used to store the sum of all the values for each key.

In line 4 A Loop which is for each is taken. It will run each time for the values inside the Iterable values. Which is coming from the second phase that is shuffle and sort phase after the mapper phase?

In line 5 Value for sum can that is “sum” is calculated and stored.

In line 7 sum cart is obtained as value to next context.

How to Execute

```
Hadoop jar topfivecategories.jar /youtubedata.txt /topfivecategories_out
```

Here ‘Hadoop’ is a command and jar specifies that we are running java type of application and topfivecategories.jar is the jar file which is created and consists of the above source code.

How to view output

```
hadoop fs-cat/topfivecategories_out/part-r-00000 | sort -n -k2 -r | head -n5
```

Here ‘Hadoop’ is a command. And df is related to the File System of Hadoop (HDFS) which is used to perform some operations on Hadoop Distributed file System. The – cat command is used to view the contents of a file and top five categories/part-r-00000 is the file where final output is stored.

This Part file is created by Text Input Format which consists of the actual final output. This Part file is created by default. Then, sort –n–k2 –r | head –n5 is the main command which shows the top 5 categories with a maximum number of videos uploaded. Also, we can mention the secondary sort instead of this command.

Sort is used to sort the data, –n stands for numerically that is sorting will be numerical, –k2 stands for the second column, –r stands for recursive operation, –n5 stands for the first 5 values after sorting.

Output





Problem Statement 2

In this problem statement, we will find the top 10 best-rated videos in YouTube.

SOURCE CODE

Now by using Mapper code, we will get the video id as key and rating as final Int value as values which will be passed to the next stage which is shuffle and sort phase and then this value is sent to the reducer phase where the aggregation of the values is performed.

Mapper Code

```

1. public class Top_Video_Rating {
2. public static class Map_1 extends
   Mapper<LongWritable, Text, Text,
3. FloatWritable> {
4. private Text top_video_name = new Text();
5. private FloatWritable top_rating = new
   FloatWritable();
6. public void map_1(LongWritable key, Text
   value, Context context )
7. throws IOException, InterruptedException {
8. String top_line = value.toString();
9.     If(top_line.length()>0) {
10.     String str_1[]=top_line.split("\t");
11.     top_video_name.set(str_1[0]);
12.     if(str_1[6].matches("\\d+.+")){
13.     float f=Float.parseFloat(str_1[6]);
14.     top_rating.set(f);
15.     }
16.     }
17. context.write(top_video_name, top_rating);
18.     }
19.     }

```

20. }

Explanation of the above code

In line 1 we are creating a class with name Top_Video_rating

In line 2 Map_1 is extended by Mapper default class. Which have the arguments keyIn, ValueIn, KeyOut, and ValueOut?

In line 4 a variable is declared which is a private Text variable top_video_name which is used to store the video name in the encrypted format.

In line 5 a variable is declared which is private FloatWritable with name top_rating which is used to store the rating of the video.

In line 6 map_1 method is overridden which is run one time for every line.

In line 8 a string variable top_line is declared which is used to store the line.

In line 9 top_line is splitting by using comma “,” delimiter and values are stored in a string array.

In line 10 an if block is declared to check whether the string array length greater than 7 or less than 7.

In line 11 video name is stored in the variable top_video_name declared in 2nd line

In line 13 numeric data is converting into float type by using type cast method

In line 14 the rating of the video is stored in the variable top_rating.

In line 17 we are writing the key and. Which will be the output of the Mapper’s Map_1 method.

Reducer Code

```

public static class Top_Reduce extends
Reducer<Text, FloatWritable, Text,
FloatWritable> {
public void Top_Reduce(Text key,
Iterable<FloatWritable> values, Context context)
throws IOException, InterruptedException
{
float top_sum = 0;
Int l=0;
for (FloatWritable top_val : values) {
l+=1;
top_sum += top_val.get();
}
Top_sum=top_sum/l;
context.write(key, new
FloatWritable(top_sum));
}
}

```

In Reducer code

line 1 a class Top_Reduce which is extended the default Reducer class with arguments KeyIn , ValueIn, same as the outputs of the mapper class and , ValueOut that is used to store final outputs of our MapReduce program.

In line 2 Reduce method is overridden which is run each time for every key. In line 4 An integer variable top_sum is

declared which is used to store the sum of all the values for each key.

In line 5 another variable is as "l" which is incremented every time. This variable is incremented as many values are there for the key.

In line 6 A Loop which is for each is taken. It will run each time for the values inside the Iterable values. Which is coming from the second phase that is shuffle and sort phase after the mapper phase?

In line 8 Value for top_sum that is "sum" is calculated and stored.

In line 10 An average of the obtained top sum is performed.

How to Execute

```
hadoop jar video_rating.jar /youtubedata.txt /videorating_out
```

Explanation for this statement is same as in Problem statement 1.

How to Execute

```
Hadoop jar topratingvideos.jar /youtubedata.txt / top rating videos _out
```

Here 'Hadoop' is a command and jar specifies that we are running java type of application and topratingvideos.jar is the jar file which is created and consists of the above source code.

How to view output

```
hadoop fs-cat/ topratingvideos_out/part-r-00000 | sort -n -k2 -r | head -n5
```

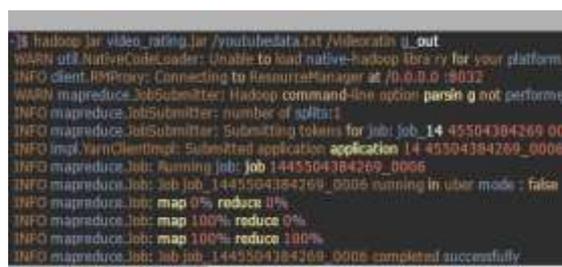
Here 'Hadoop' is a command. And df is related to the File System of Hadoop (HDFS) which is used to perform some operations on Hadoop Distributed file System. The - cat command is used to view the contents of a file and top rating videos /part-r-00000 is the file where final output is stored.

This Part file is created by Text Input Format which consists of the actual final output. This Part file is created by default. Then, sort -n -k2 -r | head -n5 is the main command which shows the top 5 categories with a maximum number of videos uploaded. Also, we can mention the secondary sort instead of this command.

Sort is used to sort the data, -n stands for numerically that is sorting will be numeric,

-k2 stands for the second column, -r stands for recursive operation, -n5 stands for the first 5 values after sorting.

Output



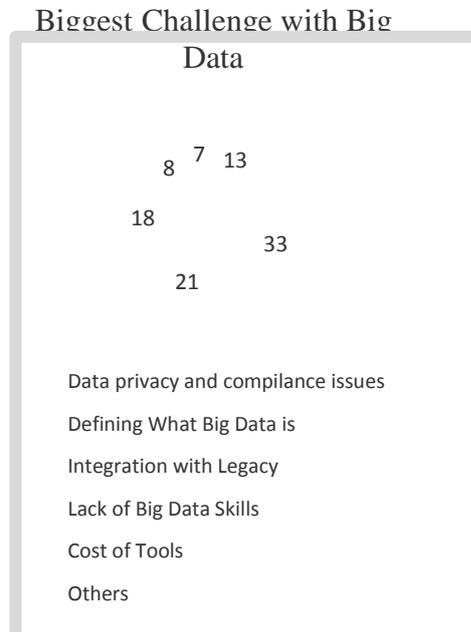
4. BIG DATA CHALLENGES

Understanding and Utilizing Big Data – The main challenge in using the Big Data approach is to understand the data that is available to process and use. It is a very challenging task to understand the data available in Information Technology Industries. This type of analysis is very important for any organization and it is performed on regular basis in order to keep pace with the market and to develop the organization [17].

Privacy, Security, and Regulatory Considerations – We all know the volume and complexity of data are becoming more complex. Keeping this in mind, it is very difficult to obtain a reliable content or data. Also, it is very difficult to prevent the data from the threats [10]. The various Security mechanism is available today. But implementing that mechanism may increase the cost of a company.

The Need for IT, Data Analyst, and Management Resources – To analyze the big data efficient human

resources with proper skill sets are required. The challenge in Big Data described in the first point that is understanding the big data. To understand and process the large volume of data, People with proper skill sets are required. Some of the Challenges are highlighted below in the pie chart.



5. CONCLUSION

A large volume of data availability of Big Data, have produced a very good method for data analysis. The availability of big data and its framework has enabled us to analyze the data at a very large scale. The main advantage of Big Data is its accuracy which leads to more confidence in decision making. If decisions are better they can result in greater performance, greater efficiency, increase productivity and also the cost effective that is can help in the cost reduction also can help in the reduction of risks. Also, there are a lot of challenges in implementing the big data approach. The biggest challenge does not seem to be technology itself as it is growing at a very high speed. But the main challenge is whether we have proper resources, proper skill sets, and proper dedicated hardware to implement this. If we have proper resources other challenges are also there like there are many legal like data privacy and integrity, cyber security which we need to resolve. With the help of Big Data, we can run applications with a large number of nodes with large number data which may be in terabytes. This allows the very fast transfer of data between nodes. It is also a very reliable framework which allows continuing operating if even any node fails. It also minimizes the risk and chance of system failure even if a good number of nodes get failed.

Nowadays, a number of companies have implemented this Big Data approach. They are working on it and at a great

pace. Companies like Facebook which have a large volume of data uses this approach.

Finally, social media sites like Facebook, Twitter and LinkedIn would not exist without the big data.

ACKNOWLEDGEMENT

I would like to thank my guide Ms. Shabnam Kumari for her indispensable ideas and continuous support, encouragement, advice and understanding me through my difficult times and keeping up my enthusiasm, encouraging me and showing great interest in my thesis work, this work could not finish without his valuable comments and inspiring guidance.

REFERENCES

- [1] H.Herodotou, H.Lim, G.Luo, N.Borisov, L.Dong, F.B.Cetin, and S. Babu. Starfish: A Self-tuning System for Big Data Analytics. In CIDR, pages 261–272, 2011.
- [2] Shweta Pandey, Vrinda Tokekar. The prominence of MapReduce in BIG DATA Processing. In Fourth International Conference on Communication Systems and Network Technologies, IEEE, pages 555-560, 2014.
- [3] Hadoop. <http://hadoop.apache.org/>.
- [4] Hadoop MapReduce Tutorial. http://hadoop.apache.org/common/docs/r0.20.2/mapred_tutorial.html
- [5] T. Nykiel, M. Potamias, C. Mishra, G. Kollios, and N. Koudas. MRShare: Sharing Across Multiple Queries in MapReduce. PVLDB,
- [6] <http://www.drdoobs.com/parallel/indexing-and-searching-on-a-hadoop-distr/226300241>
- [7] <http://www.pgs-soft.com/harnessing-big-data/>
- [8] Cohen, J., Dolan, B., Dunlap, M., Hellerstein, J.M., Welton, C.: MAD Skills: New Analysis Practices for Big Data. Proceedings of the ACM VLDB Endowment 2(2), 1481–1492 (2009).
- [9] Cuzzocrea, A., Song, I., Davis, K.C.: Analytics over Large-Scale Multidimensional Data: The Big Data Revolution! In: Proceedings of the ACM International Workshop on Data Warehousing and OLAP, pp. 101–104 (2011).
- [10] Elgendy, N.: Big Data Analytics in Support of the Decision Making Process. MSc Thesis, German University in Cairo, p. 164 (2013).
- [11] EMC: Data Science and Big Data Analytics. In: EMC Education Services, pp. 1–508 (2012).
- [12] He, Y., Lee, R., Huai, Y., Shao, Z., Jain, N., Zhang, X., Xu, Z.: RCFile: A Fast and Space-efficient Data Placement Structure in MapReduce-based Warehouse Systems. In: IEEE International Conference on Data Engineering (ICDE), pp. 1199–1208 (2011).
- [13] Herodotou, H., Lim, H., Luo, G., Borisov, N., Dong, L., Cetin, F.B., Babu, S.: Starfish: A Self-tuning System for Big Data Analytics. In: Proceedings of the Conference on Innovative Data Systems Research, pp. 261–272 (2011).
- [14] Kubick, W.R.: Big Data, Information, and Meaning. In: Clinical Trial Insights, pp. 26–28 (2012).

- [15] Lee, R., Luo, T., Huai, Y., Wang, F., He, Y., Zhang, X.: A smart: Yet Another SQL-to-MapReduce Translator. In: IEEE International Conference on Distributed Computing Systems (ICDCS), pp. 25–36 (2011).
- [16] Plattner, H., Zeier, A.: In-Memory Data Management: An Inflection Point for Enterprise Applications. Springer, Heidelberg (2011).
- [17] Russom, P.: Big Data Analytics. In: TDWI Best Practices Report, pp. 1–40 (2011).
- [18] TechAmerica: Demystifying Big Data: A Practical Guide to Transforming the Business of Government.