# An Architecture for Metadata Extractor of Big Data in Cloud Systems

Fathy E. Eassa, Hassanin Al-Barhamtoshy, Abdullah Almenbri, Osama H. Younis, Kamal Jambi

**Abstract**— In this research, we introduce an agent-based architecture for metadata extractor of big data, which is a part of a new framework for managing big data. In the architecture, many agents should be generated simultaneous to be migrated to existing big-data with different types: structured, semi-structured, and unstructured on the remote machines to collect the metadata and returned back to the user. The architecture includes the metadata storage that contains the metadata of existing big data: structured, semi-structured, and unstructured. These metadata describes the existing big data to be processed by software agents for knowledge discovery. The architecture consists of many mobile and stationary agents. The mobile agents migrate to remote machines that include big data to process and discover the required knowledge and return back to the main server. In the main server, the knowledge returned to the user by the interface agent. In this research, many techniques will be built for collecting metadata of different data such as natural language processing techniques, data mining technique, and lexical analysis technique, image processing technique, etc. The manager is built to manage big data on cloud distributed systems. In this paper, the technique of collecting metadata of tweets has been implemented and tested.

**Index Terms**— big data, mobile agents, metadata extraction, cloud, structured data, semi-structured data, unstructured data.

————————————— ◆ —————————————

## 1 INTRODUCTION

Organizations with data warehousing solutions are facing problems of data outburst. Data sets being collected and analyzed for business intelligence are growing rapidly and size of the databases used in enterprises has been growing exponentially. Sources generate data from business processes, transactions, social networking sites, web servers, sensors that exist and collected as meta-data and remains in structured, semi-structured as well as unstructured form [1][3].

Big data in cloud is one of the major issues in computing to-day, detection of global weather patterns, social phenomena, or economic changes are examples of big data analysis tasks [4]. Cloud computing, poses an important impact on industry of information technology, and research communities [5]. Many of cloud services require users to share their data like health records for data mining and analytical process, taken into consideration privacy and security [6][7]. Enterprises are looking forward for applications having large scale, web-oriented, intensive features to be accessed from diverse devices including mobile devices. Extracting meaningful information, collecting, processing, and analyzing the huge amount of data as large datasets is a challenging task. The most popular open-source map-reduce implementation framework is Apache Hadoop [2], it has been used as an alternative to store and process extremely large data sets on commodity hardware. Scalability, unstructured data, accessibility, real time analytics, fault tolerance are various challenges faced in large data management. Variations means the amount of data stored in different sectors across domains, their data types structured or unstructured examples include graphs, messages, images, audio, video, text/numeric information or data. Data types vary across enterprises [3].

In this paper, we introduce an agent-based manger for extracting metadata of existing big data to be stored in metadata storage. The stored metadata will be used for discovering knowledge and information from big data. Extracting metadata and discovering knowledge based on agent technology solves the transportation and management challenges of big data.

## 2 RELATED WORK

Extracting meaningful knowledge from big data sets is of high importance for organizations. A business has to adapt to some flexible means in terms of infrastructure and software tools. Methodologies or techniques that are being applied to big data so far include data mining grids, massive parallel processing, scalable storage systems, cloud computing platforms, distributed file systems etc. Online transaction processing (OLTP) or Online analytical processing (OLAP), with time based event or transaction time knowledge delivery is the core characteristics in Big Data Analytics. Decision makers or various businesses always look forward for exceptional methods which can process large datasets within tolerable elapsed time. Their formulation uses statistical, mathematical, algebraic and economic operations. Currently computing infrastructures are using various methods and visualization techniques to perform Big Data Analytics [8].

Apache Hadoop is an open-source software solution for scalable and reliable distributed computing. Apache Hadoop's framework contains a library that allows for handling of large datasets across clusters with hundreds of nodes, using a simple programming model. It enables applications to work with thousands of computational independent computers and for petabytes of data. Hadoop was derived from Google's MapReduce and Google File System (GFS) [2].

Hadoop Distributed File System (HDFS) [9] in distributed environment provides fault tolerance and runs on commodity hardware. Data is chopped, stored and is scattered over numerous nodes. HDFS has one master node and multiple slaves' nodes. The name node stores the metadata and data nodes store data blocks.

---

• *Authors are faculty staff in King Abdulaziz University, faculty of Computing and Information Technology, Saudi Arabia. E-mails, respectively: feassa@kau.edu.sa, hassanin@kau.edu.sa, ammali@kau.edu.sa, osama.26@live.com, kjambi@kau.edu.sa*

All of this architecture resides on commodity hardware where each node/server provides local storage and computation. To store data Hadoop distributed file system (HDFS) uses multiple nodes across networks. On top of it map-reduce framework supports and is a mean to execute jobs across nodes. HDFS has master/slave architecture. Large data is automatically splitted across nodes to be stored and retrieved within hadoop clusters.
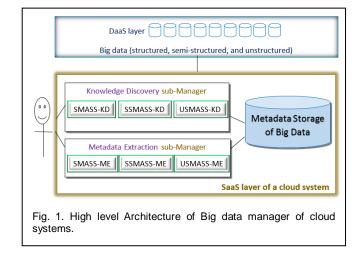
In the year of 2004, Google introduced a mapReduce programming framework for managing large data sets [10]. This framework provides ease in processing for distributed environments where data is divided into chunks and is spread over multiple clusters of thousands of nodes as map-reduce jobs. Google's model for executing large data set jobs uses a map function which in turn processes value pairs (key, index and further intermediate pairs) for data. Reduce function merges values, removes duplication for similar intermediate keys [11].

Focusing on the velocity of Big Data, a popular open-source stream processing engine (Storm) is used to perform real integration and trend detection on Twitter and Bitly streams [12]. Also, ClowdFlows platform with the real-time data streams is used to create specialized type of workflow component and a stream mining [13].

## 3 THE HIGH LEVEL ARCHITECTURE OF BIG DATA FRAMEWORK

The big data manager (shown in figure 1) consists of two sub-managers: metadata extraction sub-manager and knowledge discovery sub-manager. The metadata extraction sub-manager extracts and retrieves metadata of the big-data on the cloud machines and stores them in metadata storage. This sub-manager consists of multi-agent subsystems: multi-agent subsystem for unstructured data (USMASS-ME), multi-agent subsystem for semi-structured data (SSMASS-ME), and multi-agent subsystem for structured data (SMASS-ME). Each subsystem has many agents; the number of agents depends on the type of data that are processed by the subsystem. For example, the unstructured data includes text, videos, images, tweets, E-mails, and others. The construction of agents depends on extraction technique, type of data, and format of existing big-data.

The knowledge discovery sub-manager discovers the required knowledge or information that is needed by the user. The discovery sub-manager consists of three multi-agent subsystems: unstructured multi-agent sub-system(USMASS-KD), semi-structured multi-agent subsystem (SSMASS-KD), and structured multi-agent subsystem (SMASS-KD). Each subsystem consists of stationary and mobile agents. The task of mobile agent depends on the user query, type of big-data and the technique of discovery. For example, there is a technique for processing texts, another for processing E-mails, and so on. The discovered knowledge and information should be collected and delivered to the user.



Fig. 1. High level Architecture of Big data manager of cloud systems.

### 3.1 The Agent-based Architecture of Metadata Extraction sub-Manager

Here, we propose the agent-based architecture of the big data Metadata Extraction sub-Manager, which consists of many agents: stationary and mobile agents, see figure 2.



Fig. 2. Agent-based Metadata Extraction Manager of Big data.

Each mobile agent migrates from the big data server machine to a remote machine to collect the required data (metadata) from existing big-data. There are mobile agents for structured data, mobile agents for semi-structured data, and mobile agents for unstructured agents. The mechanism of extracting the required metadata depends on the type of the data. For example, the technique of collecting metadata from a text as unstructured data is different than collecting metadata from E-mail, Tweets, images, videos or other unstructured data. This means we have many different mobile agents for unstructured big data. Also, we have different agents for collecting metadata from semi-structured big data.

The details of some mechanisms, algorithms, and techniques of collecting metadata will be discussed here in this section. The interface agent receives the requests from the user of the big-data manager and displays the outputs. Based on the received requests, the launcher agent creates instances from different types to be migrated into remote machines in the cloud systems for extracting metadata. The agents return back to the main machine and send the extracted metadata to the Metadata Storage of Big Data (MSBD) Agent. The MSBD agent stores the returned metadata in the Metadata Storage of Big Data.

Figure 3 below shows the main activities that have been conducted extract the metadata from big data in a cloud distributed system. At the metadata server where all stationary and mobile agents are installed, many instances of mobile agents are created by the launcher agent. All instances of mobile agents migrate simultaneously into machines in the cloud system to analyze and extract the metadata of big-data.



Fig. 3. Activity Diagram of the process of Retrieving Metadata of Remote.

The algorithms and techniques of analyzing big data that will be implemented in mobile agents depend on the varieties of big-data( structured, semi-structured, and unstructured) and type of data such as free-text, images, e-mails, chats , xml, html, or databases. The natural language processing algorithms and techniques, data mining algorithms, or statistical techniques can be used for analyzing and extracting metadata. Also, the lexical analysis, data mining and parsing techniques can be used especially for analyzing and extracting the metadata from semi-structured. After the analyzing and extracting metadata activity, all mobile agents return back to the metadata server for giving retrieved metadata to the MSBD Agent to be stored in metadata storage. In this paper the technique of extracting metadata from tweets has been introduced.

## 3.2 The Metadata structure

The proposed structure of the metadata of the framework is represented as a tree as shown in figure 4. The details of the metadata will be investigated in this research.



Fig. 4. The Structure of the Metadata Storage.

## 4 IMPLEMENTATION OF TWEETS METADATA EXTRACTOR

This section presents the way of extracting the tweets' metadata from the social media (Twitter). We have used API called twitter4j (version: 3.0.3) which allows to access and retrieve the tweet's information using Java programming language. It can be used to retrieve different information related to tweet and users such as user name, display name tweet text, date of tweet …etc, and it supports the searching about tweets. Figure 5 shows the elements of the implemented extractor of metadata of Tweets.



Fig. 4. The Structure of the Metadata Storage.

The process of extracting tweets metadata is as follows:
1- *TwitterExtractor* agent calls the Twitter4j API functions for retrieving the required tweets from the Twitter.
2- *TwitterExtractor* agent applies the algorithm of extraction of tweets metadata (Algorithm 1) to:
  a) Retrieve the tweets by calling API function and passing the data in the domain table (table 1) as arguments to the called function.
  b) Extract the required metadata from the retrieved tweets.
  c) Store the extracted metadata in local data structure to be returned with the agent to the main machine. The returned metadata is stored by MSBD agent in Tweets table (Table 2) in the metadata storage of big data (figure 2)

TABLE 1
DOMAIN TABLE

| Domain Table | |
|---|---|
| Domain-No | Is a Number representing the domain number (Primary key) |
| Domain | Is a text represent the required domain (hashtag, keyword, domain, topic...etc.) |

TABLE 2
TWEETS TABLE

| Tweets Table | |
|---|---|
| tweetNo | Is a serial number of tweet |
| tweenOwner | User who written the tweet |
| screenName | Screen name for user who written the tweet |
| tweetText | Tweet content |
| tweetDate | Created date of tweet |
| domain_NO | Domain number (foreign key for relationship with the Domain Table) |

The following algorithm shows the steps of extracting metadata from tweets.

ALGORITHM 1
EXTRACTING METADATA OF TWEETS

```
Begin
    /* create object of ConfigurationBuilder */
1: builder = new ConfigurationBuilder;
    /*set the four authorizations for this object)*/
2: builder.SetOAuthAccessToken(AccessToken);
3: builder.SetOAuthAccessTokenSecret(AccessTokenSecret);
4: builder.SetOAuthConsumerKey(ConsumerKey);
5: builder.SetOAuthConsumerSecret(ConsumerSecret);
    /*create OAuth authorization object*/
6: auth=new OAuthAuthorization;
    /*assign the builder to auth*/
7: auth=builder.build();
8: tweet = new Twitter;
    /*assign auth to tweet*/
9: tweet.authorization = auth;
    /*prepare required query*/
10: queries[] = domainTableQueries;
11: for each query in queries[] do
    11.1: queryResult = tweet.search(query);
    /*store tweets info to list of Status*/
    11.2: tweets[] new Status(queryResult.getTweets());
12: end;
13: for each tweet in tweets[] do
    13.1: tweetInfo[] = tweet.userName, tweet.screenName,
                tweet.text, tweet.createdDate, tweet.domainNumber;
    /*store info to the DB*/
    13.2: tweetTable.insert(tweetInfo);
14: end;
End;
```

## 5 TESTING THE TWITTER EXTRACTOR

Test cases have been conducted to test the functionality of the extractor. The data in the domain table (Table 3) used as arguments for retrieving tweets to be analyzed to extract the metadata. The results returned based on these queries are stored in different tables shown in the Tables 4, 5, 6.

TABLE 3
DOMAINTABLE



TABLE 4
THE RETRIEVED RESULTS FOR DOMAIN "CRISIS IN SYRIA"



TABLE 5
THE RETRIEVED RESULTS FOR DOMAIN "CRISIS IN EGYPT"

TABLE 6
THE RETRIEVED RESULTS FOR DOMAIN "SAUDI ARABIA"



In addition to the above tests, extracting metadata of tweets written in Arabic language has been conducted as shown in Table 7.

TABLE 7
THE RETRIEVED RESULTS FOR ARABIC DOMAIN (مستشفى_العرضي , البرد في سوريا)



# 6  CONCLUSION

In this research, an agent based sub-system for collecting the metadata of big data has been introduced. The sub-system is part of a big data manager that solves two big data challenges: big data transportation, and big data management. The transportation challenge will be solved by implementing part II of our manager. Part I that has been introduced collects the metadata of existing big data. The metadata of big data is different from type to other. Therefore, the design and techniques of mobile agents for collecting the metadata are different. In this paper we introduced the algorithm and technique of collecting metadata of tweets. The algorithm has been implemented as agent called Twitter Extractor Agent. . The agent has been tested and the collected metadata stored in the metadata storage. Our big data manager is scalable because many instances of agents can be created at the same time to be migrated to many machines. Therefore, the manager can man-

age any size of big data stored on any number of storage media and machines.

In the future, many extractors will be built to enhance the implementation of the manager. For example we will built extractor for E-mails, extractor for text, extractor for XML, and so on.

## REFERENCES

[1] Impetus white paper, March 2011, "Planning Hadoop/NoSQL Projects for 2011" by Technologies, Available at: http://www.techrepublic.com/whitepapers/planninghadoopnosql-projects-for-2011/2923717.

[2] Apache Hadoop. [Online] Available at: http://wiki.apache.org/hadoop. [Accessed Jan 2014].

[3] A. Jacobs. (2009, Jul.) The pathologies of big data. [Online]. Available: http://queue.acm.org/detail.cfm?id=1563874 [Accessed Jan 2014].

[4] Chamikara Jayalath, Julian Stephen, and Patrick Eugster, (2014). From the Cloud to the Atmosphere: Running MapReduce across Data Centers, IEEE TRANSACTIONS ON COMPUTERS, VOL. 63, NO. 1, JANUARY 2014, pp. 74 - 87.

[5] Xuyun Zhang, Laurence T. Yang, Senior Member, Chang Liu, and Jinjun Chen, (2014). A Scalable Two-Phase Top-Down Specialization Approach for Data Anonymization Using MapReduce on Cloud, IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 25, NO. 2, FEBRUARY 2014, pp. 363 - 373.

[6] D. Zissis and D. Lekkas, "Addressing Cloud Computing Security Issues," Future Generation Computer Systems, vol. 28, no. 3, pp. 583-592, 2011.

[7] L. Hsiao-Ying and W.G. Tzeng, "A Secure Erasure Code-Based Cloud Storage System with Secure Data Forwarding," IEEE Trans. Parallel and Distributed Systems, vol. 23, no. 6, pp.

[8] McKinsey Global Institute, 2011, Big Data: The next frontier for innovation, competition, and productivity, Available: www.mckinsey.com/~/media/McKinsey/dotcom/Insights%20and%20pubs/MGI/Research/Technology%20and%20Innovation/Big%20Data/MGI_big_data_full_report.ashx, Aug 2012.

[9] Apache Software Foundation. [Online] Official apache Hadoop website: http://hadoop.apache.org/, [Accessed January 2014].

[10] The Hadoop Architecture and Design. [Online], Available: http://hadoop.apache.org/common/docs/r0.16.4/hdfs_design.html, [Accessed Jan 2014].

[11] Hung-Chih Yang, Ali Dasdan, Ruey-Lung Hsiao, and D. Stott Parker from Yahoo and UCLA, "Map-Reduce-Merge: Simplified Data Processing on Large Clusters", paper published in Proc. of ACM SIGMOD, pp. 1029–1040, 2007.

[12] Thibaud Chardonnens, Benoit Perroud, (2013). Big Data Analytics on High Velocity Streams: A Case Study, 2013 IEEE International Conference on Big Data, pp. 784 – 787.

[13] Kranjc, Janez ; Podpecan, Vid ; Lavrac, Nada, (2013). Real-time data analysis in ClowdFlows, 2013 IEEE International Conference on Big Data, pp. 15 – 22.