



An Efficient Mapreduce-based System to Find Userlikeness on Social Networks

By D. Ravikiran & Dr. S.V.N Srinivasu

Acharya Nagarjuna University, India

Abstract- Day to day Social network information growth pursues an exponential pattern, and Present DB management systems cannot manage efficiently such a huge volume of data. It is essential to employ a “big data” solution for Social network problems. One of the most important problems in Social network is finding User likeness (ULi). Current methods for finding ULi are not flexible and do not sustain all data sources, nor can they accomplish user necessities for a query tool. In this paper, we propose a reliable and data available method to solve ULi problems over MapReduce design. RiDaULi supports storage and retrieval of all kinds of data sources in an appropriate manner. The dynamic nature of the proposed method helps users to define conditions on all entered fields. Our assessment shows that we can use this method as high confidence in less execution time.

Keywords: social networking, userlikeness, mapreduce, mapper.

GJCST-C Classification : K.6.3 D.2.12



AN EFFICIENT MAPREDUCE BASED SYSTEM OF FINDING USER LIKENESS ON SOCIAL NETWORKS

Strictly as per the compliance and regulations of:



An Efficient Mapreduce-based System to Find Userlikeness on Social Networks

D. Ravikiran^α & Dr. S.V.N Srinivasu^ο

Abstract- Day to day Social network information growth pursues an exponential pattern, and Present DB management systems cannot manage efficiently such a huge volume of data. It is essential to employ a “big data” solution for Social network problems. One of the most important problems in Social network is finding User likeness (ULi). Current methods for finding ULi are not flexible and do not sustain all data sources, nor can they accomplish user necessities for a query tool. In this paper, we propose a reliable and data available method to solve ULi problems over MapReduce design. RiDaULi supports storage and retrieval of all kinds of data sources in an appropriate manner. The dynamic nature of the proposed method helps users to define conditions on all entered fields. Our assessment shows that we can use this method as high confidence in less execution time.

Keyword: social networking, userlikeness, mapreduce, mapper.

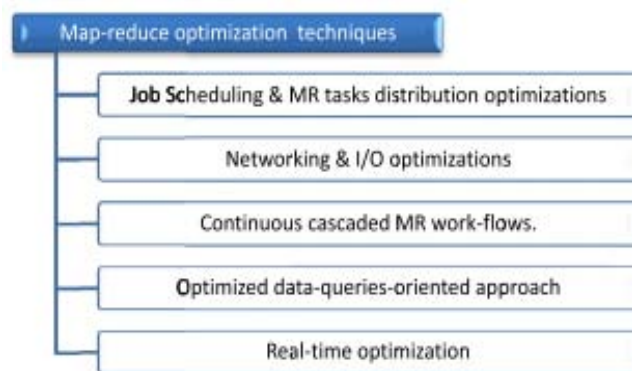
1. INTRODUCTION

Now a days, with huge volume of user data contraction, common or frequent database management systems cannot effectively sustain data management and analysis in many fields, including meteorology, scientific instruments, social networks, and medical networks. In these and other fields we need a pattern shift to address our problems. Capturing, storing and retrieving information in a timely manner are vital issues in these systems. It is necessary to have available and reliable solutions for these kinds of problems because the prevalent single-node and parallel approaches are far from offering a timely solution. On the other hand, reliable and available resolutions have their own troubles, in particular network bottlenecks, low performance of hardware nodes, and necessities for other nodes' information. Social Network is one of the fields that need reliable and data available solutions, because current solutions cannot properly solve this area's problems. One of the most important problems in this area is identifying user's likeness, or ULi, defined as the rate of likeness between two or more users in terms of their like, interests, personal information, etc. The goal in ULi is to identify those Users who have the greatest amount of information in common in order to use their Preferences or recommendations for new users.

We have two main issues in ULi: the huge amount of information per users; and the fact that most

of this data is nonstructured, lacking a predefined record structure that is common among all users. A large number of fields per users may add complexity to ULi problems as well. Given these characteristics, we have to use so-called “big data” solutions. One of the methods which can be used for reliable and data available solutions for big data is MapReduce. MapReduce is used to solve Social Network problems. But MapReduce and other data available solutions have problems such as data locality, network bottlenecks, hardware inefficiency etc. In this paper, we propose RiDaULi, a reliable and data available method for investigating user's likeness. In this method, a MapReduce-based method is used to solve ULi problems. Unlike other approaches, we do not use structured or semi-structured methods for user's information storage. RiDaULi can use different data sources with different data items. Even the same data source can have different data items for two users. Rather, RiDaULi uses a dynamic method to store user's information which can be easily dispersed over hardware nodes. In the proposed method hardware nodes can execute their tasks simultaneously, and none of the nodes needs information from other nodes which is the main problem of MapReduce-based methods. The structure of this paper is as follows. Section 2 investigates some preliminaries concerning MapReduce and Social Network problems. In Section 3, ULi-related literature is discussed. Section 4 focuses on the proposed method. Section 5 presents the evaluation of the proposed method. Section 6 provides the conclusion.

Fig 1: MapReduce Optimization



Author α: Acharya Nagarjuna University.
e-mails: ravikiran2005@gmail.com, drsvnsrinivasu@gmail.com

II. GROUND WORK

In this part, both MapReduce and the relationship between Social Network and big data are explained.

a) MapReduce

In this section, the literature related to MapReduce design is discussed, a decomposable algorithm, partitionable data, and sufficient small data partition are the main characteristics required for effective use of MapReduce. In [23], classic MapReduce was optimized to decrease the data transformation load. In the method described in [23], a shared area for information was considered. This type of design is suitable for solving problems, such as k-nn and top k queries. MPI (Message passing interface) was used for message passing in a MapReduce structure. The goal of that paper was to decrease the amount of data transferred in the MapReduce network. A method was developed for tackling workloads in hierarchical MapReduce architectures. Hadoop and uses a deduplication-based snapshot differential algorithm (D-SD) and update propagation. Haloop is another type of MapReduce structure suitable for iterative problems. iMapreduce also supports iterative processes. In [20], HDFS (Hadoop file system) was substituted with a concurrency optimized data storage layer based on the BlobSeer data management service. In [22], a model was presented to estimate I/O behavior of MapReduce applications. In [21], optimization over MapReduce

structure was divided into five groups. Fig. 1 shows these groups

b) Social Network and big data

In this section, Social Network and its relation to big data are investigated. These days, users' information is generated at an exponential rate. This information has different formats and standards. According to [19], there are various standard data sources, As shown in Fig. 2, huge Volume of information is generated in Various formats with high Velocity; therefore, we have three Vs of Big data in Social Network networks. With ULi there is an additional challenge, namely Veracity, meaning that for many users we typically have doubtful or uncertain information. Social Network problems visible all of the V's, and therefore it is inevitable that we will use big data solutions to solve them but, according to [19], existing big data technologies do not effectively deal with the full spectrum of Social Network problems, so it is necessary to customize them for our purposes. According to high volume of information in Social Network big data is necessary for data analysis .Also costs are reduced by using big data analytics in Social Network. In a users-centered framework is proposed that Can personalize Social Network with a big data driven approach. In [35] big data is used to solve problems like the selection of appropriate recommendation paths or improvement of Social Network systems. AITION [37] proposed a reliable knowledge data discovery platform for big data Social Network.



Fig. 2 : Standard data sources in Social Network

III. LITERATURE ON ULI

In this section, literature specifically concerned with ULI is investigated. According to [1], finding ULI solutions can be divided into two parts. Fig. 3 shows this categorization. The first category is solutions that identify ULI relationships by machine learning algorithms [3–5]. These types of solutions are offline and they require a long time for the machine learning to take place. Also there are data mining methods which work on streaming data and they can be considered as online data mining methods. These methods can only work on a part of data. In other word they have methods like sliding window, sampling, synopsis etc. over stream data; therefore, this method is not appropriate for ULI problem because we need to analyze all data items [40]. The second category uses information retrieval techniques. Some techniques use simple search [6,7]; however, searching over limited keywords within a predefined structure may have severe limitations. Another information retrieval solution involves Using Entity Relationship Graphs (ERG) to investigate similarities between de- fined entities [8,9]. These types of solutions

are expensive, and some are not online [8,9]. Some methods try to improve the ERG solution by unified search [10,11]. In [2] MapReduce is used to solve the problem. They tried to reduce algorithm execution time by distributing computation on hardware nodes. PARAMO [36] is a method which uses MapReduce to develop a predictive modeling platform in the Social Network analytics domain. Some methods used LSH [39] (Locality-Sensitive Hashing) for finding similarities [31]. In [31] LSH and MapReduce are used to extract user's likeness. LSH is not suitable for ULI problem because it works with predefined data structure and with ever changing data sources accuracy will reduced dramatically. According to our investigation, none of the above-mentioned methods are fully effective for solving ULI problems, because of the following considerations: ULI requires a dynamic structure to store users' information. Different users have different data items, and thus require a structure which can store data with different standards and different data formats with no default assumptions.

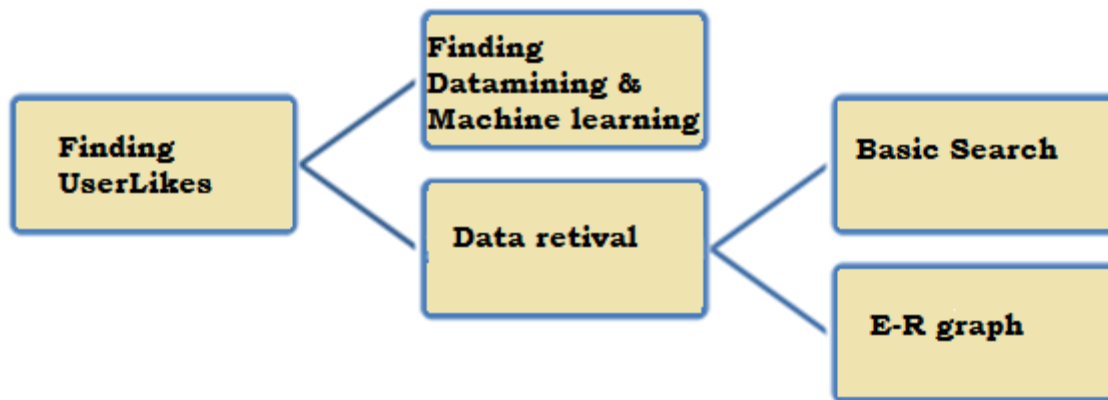


Fig. 3 : Finding user likes

- In the ULI data retrieval phase, the proposed method has to accept all types of input data items and be able to dynamically create queries over all users' data fields.
- ULI implementation time is very important; the method has to implement in a appropriate manner and with high precision. Offline and long-time query execution is not satisfactory.
- Given the huge volume of data generation, distributed solutions are necessary. In this paper we introduce RiDaULI, a reliable and data available method that uses dynamic data structure to store users' data items from data sources with different formats. It can also retrieve data items by dynamic query generation. In this connection our system achieves reliable and data available architecture of

RiDaULI, acceptable query execution time is achieved. To the best of our knowledge, RiDaULI is unique in being able to offer a solution to the ULI problem.

IV. PROPOSED METHOD

With our proposed method we illustrated RiDaULI is a reliable and data available method which is based on MapReduce. In this method, users' input data is converted to a integrated format as explained below. This adaptation has two main primitive advantages. First, varying in input data does not affect the RiDaULI format; therefore, we can allow any data format without any changes in our format. Second, this format is suitable for MapReduce architecture and helps us to dispense data over nodes. Moreover, each node can do its tasks without the need for other nodes' information.

Because of these advantages, we can easily solve ULI various formats can be stored, and efficiency can be problems over distributed nodes. Users' records in achieved by autonomous calculations.

Table 1: RiDaULi data source

Source ID	Source name
1	Facebook
2	Twitter
3	Linkedin
....

Table 2 : Input data items

Id	Name	age	Gender	habits	Likes1	Likes2
1211	sai	20	Male	Reading books	Spiritual	fiction
1212	ram	40	Male	Watching Movies	Action	comedy
1213	seetha	35	Female	Listening Music	melody	devotional

Table 3 : RiDaULicolumn

Column Id	Column name	Data Source ID
1	ID	1
2	name	1
3	age	1

Table 4 : Data(fact)

Column Id	Row Id	Value
1	1211	sai
2	1211	20
3	1211	male

Table 1 shows RiDaULiDataSource structure. In this table data source names and ID are stored.

Suppose that we are working with the information in Table 2 from Facebook data source. If we define the columns as in Table 3 (RiDaULi Column), the Table 2 data items can be converted into Table 4 (RiDaULi Fact). The data format in Table 4 has several advantages:

- Dynamic columns definition
- Completion of all fields is not necessary
- Unified data format
- Data storage size reduction

The proposed data format is suitable for the MapReduce structure, and allows us to execute queries simultaneously on different nodes. There are several steps to Using RiDaULi:

- ETL (Extract/Transform/Load): First, information from different data sources is gathered, and the metadata table (like Table 3) and data table (like Table 4) are created.

GetColumnID function retrieves ColumnID of a specific field from the RiDaULiColumn table. Input parameters are DataSourceID and ColumnName.

Also to identify equal fields on different data sources it is necessary to have the RiDaULiEqual table. Table 5 shows RiDaULiEqual.

a) Data allocation

Because of the unified data format of RiDaULi, data can be distributed over different nodes. Processing power and memory of each hardware node can be important factors to allocate data items to each node.

b) Query execution

To execute queries over MapReduce architecture, the queries first have to be converted to an appropriate format for RiDaULi. Then each converted query is sent to the nodes separately for execution, and the RowIDs of the results are returned. Finally, the extracted RowIDs are sent to the Phase 2 Mappers, and users' information is retrieved.

As shown in Fig. 4, each Phase 1 Mapper sends its results as triples. In the Phase 1 Reducer, aggregation is done on Score based on RowID, and the final Score per RowID is calculated. In the Phase 2 Mapper, other fields with corresponding RowIDs are extracted. The resulting formats of Phase 2 Mappers areas. In Phase 2 Reducer, results of Phase 2 Mappers are aggregated. Also, Phase 1 Reducer results are sent directly to the likeness Ranker, which sorts RowIDs according to their scores; then, when a RowID is selected by the user, other related information is extracted.

And ...". First all ColumnIDs are extracted from the RiDaULiColumn table. Then all rows that are equal to extracted ColumnIDs are retrieved from the RiDaULiFact table. Emit function execute queries and put results into the specified table on the specified server. If the specified table does not exist it creates a table with the specified name. For the Score calculation, many algorithms can be used. Here we use a simple algorithm, in which input users data items are compared with the same data items of existing users. If the data item value of the existing users is exactly equal to the input user's data item value, then its Score is equal to

two. Otherwise, if the user's data item value is partially similar to an existing user's data item value, then the Score is equal to one. If there is no likeness between the input data item value and the existing data item values then the Score is equal to zero. In the data sources there are many misspellings, imprecise terms, colloquial terms, etc. To solve these problems we use metadata to create associations between columns. In the Query builder phase, we can define column groups which contain the main term together with its colloquial terms, imprecise terms and prevalent misspellings. When an input column is used in a query, all other

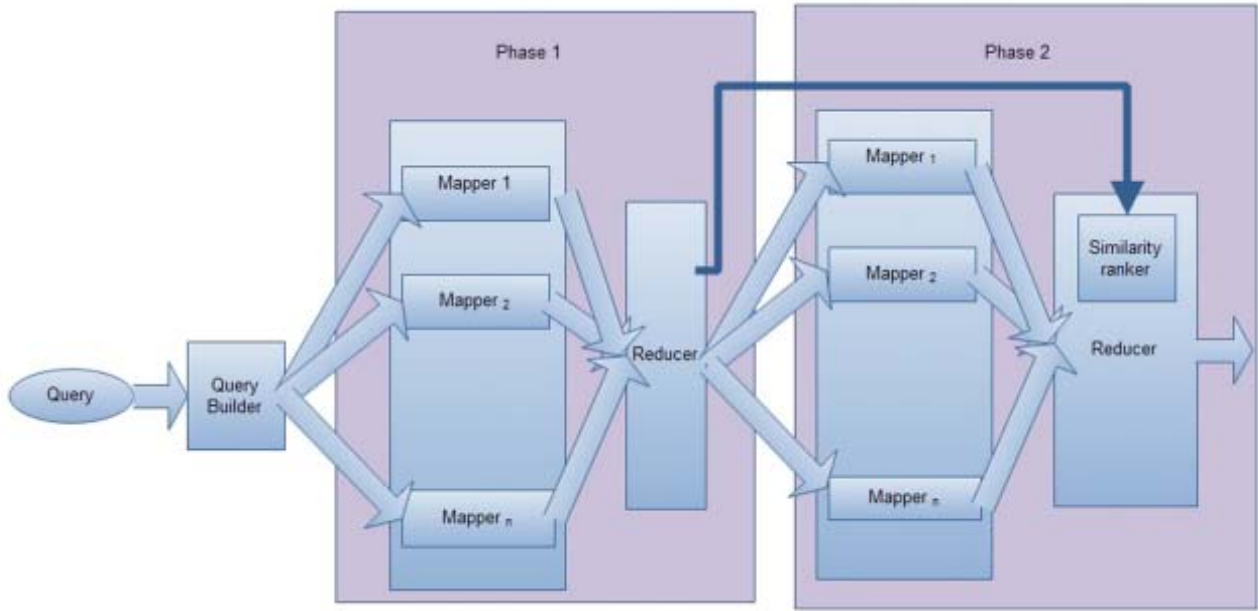


Fig. 4 : RiDaULi Process to execute query

Group members are considered and their related information is gathered. If there is a bottleneck in

the Reducer phase, we remove these via combiners. Fig. 5 shows the RiDaULi architecture with combiners.

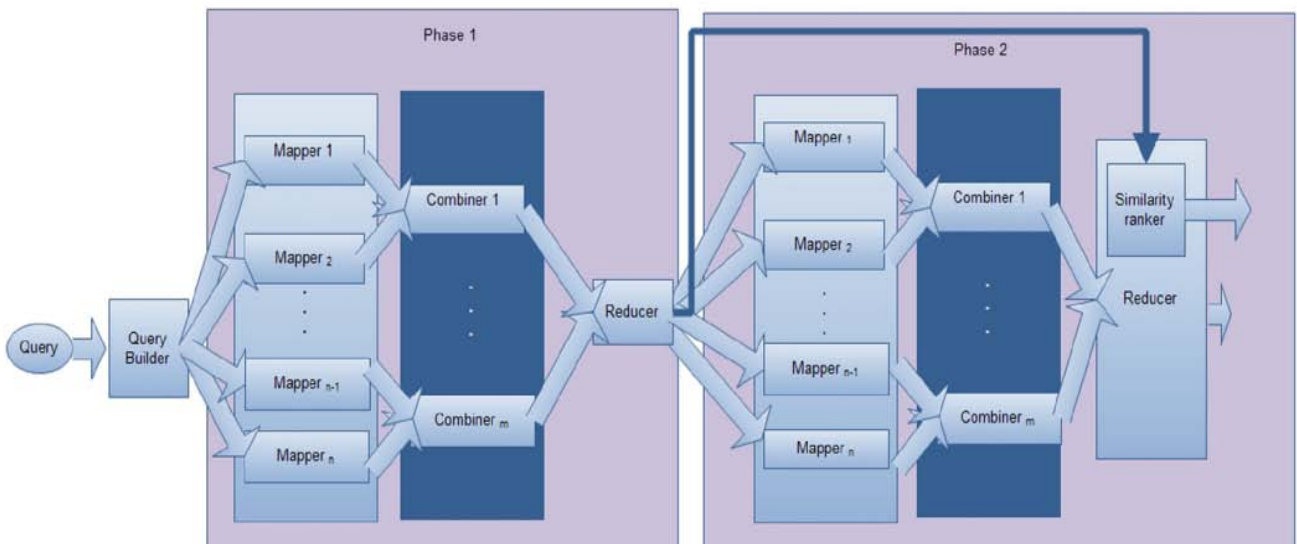


Fig. 5 : RiDaULi architecture with combiner

V. EVALUATION

In this section we evaluate RiDaULi from two views. First the execution time of the proposed method is evaluated, and second the accuracy of RiDaULi is calculated. As per illustration we producing sample Expected results.

a) Execution time

In this We used data from different Social Network systems, which in turn have different standards for storing data, by Using RiDaULi, we found that we

could easily achieve the required results on a reliable and data available structure. As shown in Fig. 4, twenty-one servers were used in Phase 1 and twenty-one for Phase 2. For thirty seven different queries we achieved an average time of 9.42 seconds. As shown in Fig. 5, we then added five combiner servers with the same specifications to each of the two phases, for a total of 52 servers. The average execution time for thirty seven queries improved about 60%, decreasing to 5.65 seconds.

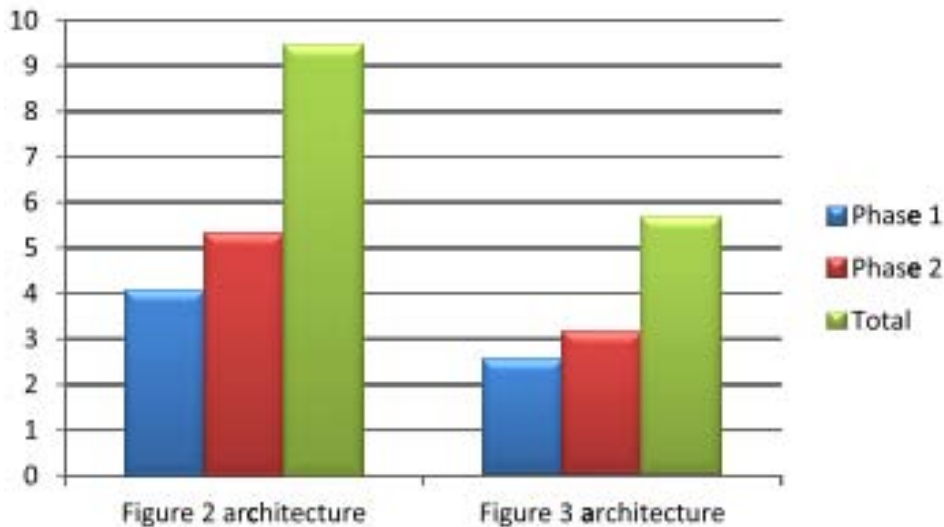


Fig. 6 : RiDaULi architecture with two combiners

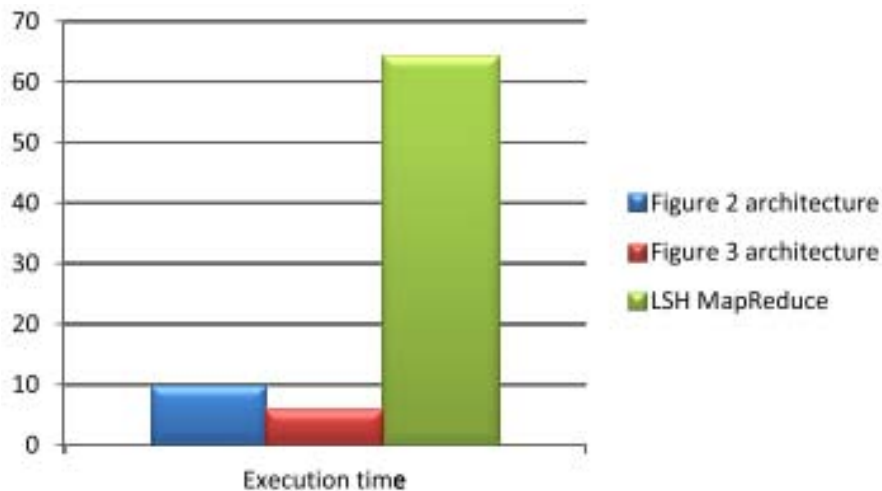


Fig 7 : Total Execution Time

Fig. 6 shows a comparison between the two phases of RiDaULi shown in the architectures of Figs. 4 and 5. Also we used the LSH algorithm over MapReduce for evaluation. 52 servers with the Table 7 specification were used. For thirty seven different queries we achieved an average time of 63.11 seconds. Fig. 7 shows the results.

VI. CONCLUSION

In this paper, we propose RiDaULi, a reliable and data available method to solve user likeness (ULi) problems over Social network. Previously, the standard methods were based on Machine Learning (ML) or Information Retrieval (IR). ML methods need a long time to execute, and are offline. Standard IR methods have

many limitations for information storing and query processing; they support only a basic user interface, and limit the kinds of queries that can be built. Online data mining methods have good performance with predefined data sources and are not suitable for dynamic data sources. Also there are some methods like LSH that can properly work over distributed environments but their performances are decreased when there are many changes in input data sources. RiDaULi is an IR method which supports different data formats. All of these formats can be retrieved by data unification. In this method all fields need not be completed, and for each user only the existing fields are entered. This feature allows for data storage size to be considerably reduced. Our evaluation shows that RiDaULi can solve ULi problems effectively. Because of the reliable and data available nature of RiDaULi, it can utilize hardware effectively in order to solve problems involving huge amounts of data.

REFERENCES REFERENCES REFERENCIAS

1. S. Fortunato, BCommunity detection in graphs,[Phys. Rep., vol. 486, no. 3–5, pp. 75–174, 2010.
2. M. E. J. Newman, BDetecting community structure in networks,[Eur. Phys. J. BCondens. Matter Complex Syst., vol. 38, no. 2, pp. 321–330, Mar. 2004.
3. Facebook Data Tracker. [Online]. Available: <http://www.checkfacebook.com>
4. M. E. J. Newman and M. Girvan, BFinding and evaluating community structure in networks,[Phys. Rev. E, vol. 69, no. 2, pp. 026113-1–026113-15, Feb. 2004.
5. M. E. J. Newman, BModularity and community structure in networks,[Proc. Nat. Acad. Sci. USA, vol. 103, no. 23, pp. 8577–8582, Jun. 2006.
6. W. W. Zachary, BAn information flow model for conflict and fission in small groups,[J. Anthropol. Res., vol. 33, no. 4, pp. 452–473, 1977.
7. M. Girvan and M. E. J. Newman, BCommunity structure in social and biological networks,[Proc. Nat. Acad. Sci. USA, vol. 99, no. 12, pp. 7821–7826, Jun. 2002.
8. J. Dean and S. Ghemawat, BMapReduce: Simplified data processing on large clusters,[Commun. ACM, vol. 51, no. 1, pp. 107–113, Jan. 2008.
9. J. Lin and M. Schatz, BDesign patterns for efficient graph algorithms in MapReduce,[in Proc. ACM 8th Workshop Mining Learn. Graphs, 2010, pp. 78–85.
10. S. Brin and L. Page, BThe anatomy of a large-scale hypertextual Web search engine[1],[Comput. Netw. ISDN Syst, vol. 30, no. 1–7, pp. 107–117, Apr. 1998.
11. A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee, BMeasurement and analysis of online social networks,[in Proc. 7th ACM SIGCOMM Conf. Internet Meas., San Diego, CA, USA, Oct. 2007, pp. 29–42.
12. C. Wilson, B. Boe, A. Sala, K. P. N. Puttaswamy, and B. Y. Zhao, BUser interactions in social networks and their implications,[in Proc. 4th ACM Eur. Conf. Comput. Syst., Nuremberg, Germany, Mar. 2009, pp. 205–218.
13. The OSN Data Set. [Online]. Available: <http://current.cs.ucsb.edu/facebook/index.html>
14. S. E. Schaeffer, BGraph clustering,[Comput. Sci. Rev., vol. 1, no. 1, pp. 27–64, Aug. 2007.
15. W. Xue, J. Shi, and B. Yang, BX-RIME: Cloud-based large scale social network analysis,[in Proc. IEEE Int. Conf. Services Comput., 2010, pp. 506–513.
16. B. W. Kernighan and S. Lin, BAn efficient heuristic procedure for partitioning graphs,[Bell Syst. Tech. J., vol. 49, no. 1, pp. 291–307, 1970.
17. J. R. Tyler, D. M. Wilkinson, and B. A. Huberman, BE-mail as spectroscopy: Automated discovery of community structure within organizations,[Inform. Soc., vol. 21, no. 2, pp. 143–153, 2005.
18. M. J. Rattigan, M. Maier, and D. Jensen, BUsing structure indices for efficient approximation of network properties,[in Proc. 12th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2006, pp. 357–366.
19. U. Brandes, D. Delling, M. Gaertler, R. Goerke, M. Hoefer, Z. Nikoloski, and D. Wagner, BOn modularity clustering,[IEEE Trans. Knowl. Data Eng., vol. 20, no. 2, pp. 172–188, Feb. 2008.
20. N. P. Nguyen, T. N. Dinh, Y. Xuan, and M. T. Thai, BAdaptive algorithms for detecting community structure in dynamic social networks,[in Proc. IEEE INFOCOM, 2011, pp. 2282–2290.



This page is intentionally left blank

GLOBAL JOURNALS INC. (US) GUIDELINES HANDBOOK 2015

WWW.GLOBALJOURNALS.ORG