# Influential Factors for Cost Minimization in Big Data Processing

**Neeraj Saxena[1] Chethana R Murthy[2]**

[2]Assistant Professor

[1,2]Department of Information Science and Engineering

[1,2]RV College of Engineering, Bangalore-560059

*Abstract—* Practical system for huge information preparing and investigation is a structure which gives cost minimization of enormous information handling between the datacenters. These datacenters are Geo-dispersed consequently puts an overwhelming weight on correspondence, stockpiling, and reckoning, which henceforth give vast measure of operational expense to the suppliers of server farms. Accordingly, three variables, which profoundly impact the operational consumption of server farms are data resizing, task assignment, task processing. We formulate the algorithms for cost optimization considering all three variables.

*Key words:* Vawt, Magnate, Magnetic Levitation, Wind Turbine, Blade hub

## I. INTRODUCTION

In 2015, 71% of overall server farm equipment spending will originate from the enormous information handling; this is anticipated by Gartner. e.g., Google's 13 data centers over 8 countries in 4 continents [1]. Procedure of Data focus resizing (DCR) information region may bring about a misuse of assets. Less mainstream information may stay unmoving and the low asset utility causes more servers to be enacted and brings about higher working Expenditure and big data have high demand on computation [2]. A few Links in systems differ on the transmission rates and expenses as per their extraordinary features. If the directing technique among server farms fizzles then it is unavoidable to download the data from a remote server. For this situation, transmission expense is relying upon directing strategy. The Quality-of-Service (QoS) of huge information assignments has not been considered in existing work.

We are roused and impacted to tackle the expense minimization issue through a joint the expense minimization issue through a joint form streamlining of these three components for huge information benefits in geo-dispersed server farms.

In order to depict the errand fulfillment time while considering both information transmission and reckoning, we propose a markov chain which is two dimensional and infer the normal undertaking finish time which is in shut structure which give the time effectiveness. At that point issue is displayed as a mixed-integer non-linear programming (MINLP) and we propose an effective answer for linearize it so we can get high measure of expense decrease. The procedures which tackle the expense minimization are Loading and Preprocessing of Data, Identification of data centers, Task processing and Task Assignment.

Data center resizing (DCR) has been proposed to decrease the processing cost by changing the quantity of actuated servers through undertaking situation [3]. In light of DCR, a few studies have investigated the topographical circulation nature of information focuses and power value heterogeneity to bring down the power cost [4]–[6].

Enormous information administration structures, e.g.,[7], embody a disseminated record framework underneath, which conveys information lumps and their imitations over the information places for fine-grained burden adjusting and high parallel information access.

In spite of the fact that the above arrangements have acquired some positive results, they are a long way from accomplishing the cost efficient huge information handling due to the accompanying shortcomings. For instance, most calculation asset of a server with less prominent information may stay unmoving. So Task Placement and data resizing both will play a important role in cost minimization of big data processing. The data reside inside the data center where it to be processed for omitting remote data loading [7][8].
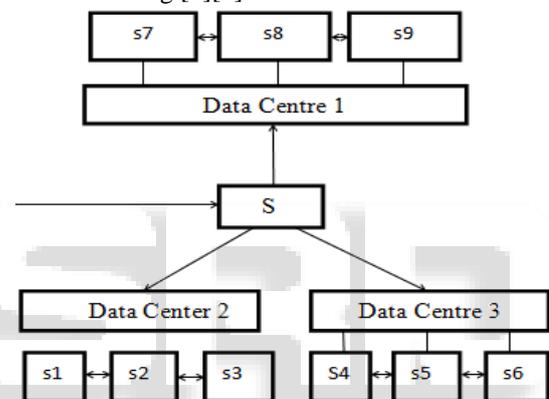


Fig. 1: Topology of server and data centers
Where, S: Scheduler   s: Servers

The above figure 1 represents the process involved in configuring the data into the processing framework and datacenter assignment.

## II. UPLOADING AND DATA PREPROCESSING

The database load utility peruses client gave information and a table stores it. There is an information document gave by client that contains information .There is four unique arrangements of records are bolstered by database load utility. Need is there to characterize a table before you can perform information stacking. The database and the tables to load probably been as of now made before stacking the information. There is an accessibility of numerous utility projects which can assemble databases and with the SQL CREATE TABLE articulations we characterize the client table. Henceforth when there is a start of stacking process, the database framework manufactures essential key records for each of the tables which have an essential key. Information ought to be preprocessed for the end invalid a quality which is available in the info dataset gave by client. In the information mining process the preprocessing of information is a vital step. The expression "waste in, rubbish out" is especially relevant to information mining and machine learning ventures.

## III. SELECTION OF DATA CENTERS

The Data Center ought to be chosen by and stockpiling limit of servers dwell in the data centers. Distinguishing proof of Data Center is vital matter for minimizing operational consumption of servers dwell in the every information centers. the links in networks vary on the transmission rates and costs according to their unique features [9].

The no. of activated and idle servers plays a vital role in selecting the data centers as the idle servers can be selected out of the data centers pool so that the cost of power Can be utilized in a effective way it can be possible. Further routing strategy matters on the transmission cost. As indicated by Jin et al. [10]

## IV. ASSIGNMENT OF TASK

Further expanding the quantity of servers won't influence the circulations of tasks. Task ought to be relegated to server farm where numbers of initiated servers are optimal. Task is profoundly impact the operational consumption of information center. Task is doled out to server farm as per closest server farm for adequately handling of data. Each information lump has a stockpiling necessity and will be needed by enormous information utilization

## V. DATA LOADING

A Data Placement on the servers and the measure of burden limit appointed to every document duplicates in order to minimize the correspondence expense while guaranteeing the client experience. Propose a mechanized data situation component Volley data focus limit limits, data between conditions, and so forth. Cloud administrations make utilization of Volley by submitting logs of datacenter for Solicitations. Volley investigates the logs utilizing an iterative enhancement algorithm taking into account data access examples and customer areas, and yields relocation recommendations back to the cloud administration. Develop Min Copy sets, a data replication situation conspire that decouples data appropriation and replication to enhance the data strength properties in dispersed data centers. As of late, Jin et propose a joint enhancement conspire that all the while upgrades virtual machine (VM) position and system stream steering to amplify vitality reserve

## VI. TASK PROCESSING

The high computational server ought not handling the low populace of information chunk. Because it expands the operational use of server, wastage of capacity and transmission cost. The populace of information is handled rely on the computational limit of servers live in the data center.
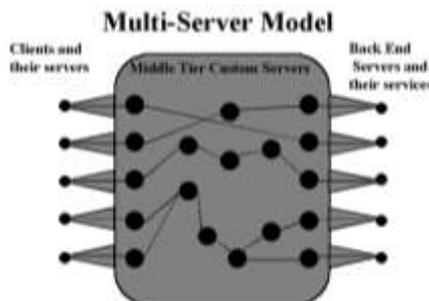


Fig. 2: Multi-Server Model
With the hadoop big data processing platform the task of processing is performed and the cost of processing for data size and arrival rate of task is calculated for analysis.

The above figure 2 represents the multi-server model of the cost minimization framework in which multiple servers can access the information from back end as well as from the front end. The middle tier servers' acts as the proxy servers for the system and these will help in circulating the information from the one end to other end with the help of the client server architecture. Architecture for the communication between servers is multi- tier client server service.

Processing of task is a major step in cost evaluation .The processing of data goes through several data centers .The servers pass the data in form of message to other servers these are known as the middle tier servers .The End to End communication between the server is done through the communication links available in the model .So other architecture which gives the complete showcase of the communication between the servers are multi tier client server model . We can depict the flow of data from the figure from both academia and industry [5], [11]–[13].
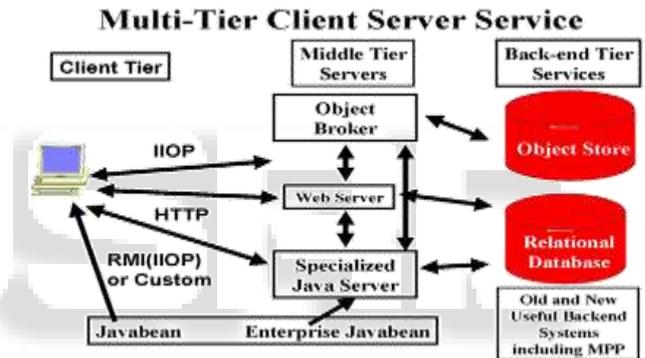


Fig. 3: Multi-Tier Client Server Service
The above figure 3 shows the communication mechanism between the client and the servers through the middle tier servers. Protocols such as IIOP, HTTP and RMI are used for message passing.

## VII. AN FORMULATION OF LINEARIZATION

| Constants | |
|---|---|
| $J_i$ | The set of servers in data center $i$ |
| $m_i$ | The switch in data center $i$ |
| $w^{(u,v)}$ | The weight of link $(u, v)$ |
| $\phi_k$ | The size of chunk $k$ |
| $\lambda_k$ | The task arrival rate for data chunk $k$ |
| $P$ | The number of data chunk replicas |
| $D$ | The maximum expected response time |
| $P_j$ | The power consumption of server $j$ |
| $\gamma^{(u,v)}$ | The transmission rate of link $(u, v)$ |
| Variables | |
| $x_j$ | A binary variable indicating if server $j$ is activated or not |
| $y_{jk}$ | A binary variable indicating if chunk $k$ is placed on server $j$ or not |
| $z_{jk}^{(u,v)}$ | A binary variable indicating if link $(u, v)$ is used for flow for chunk $k$ on server $j$ |
| $\lambda_{jk}$ | The request rate for chunk $k$ on server $j$ |
| $\theta_{jk}$ | The CPU usage of chunk $k$ on server $j$ |
| $\mu_{jk}$ | The CPU processing rate of chunk $k$ on server $j$ |
| $f_{jk}^{(u,v)}$ | The flow for chunk $k$ destined to server $j$ through link $(u, v)$ |

Table 1: Notations

The table 1 shows the constants and variables used in the formulation of the cost needs to be calculated .These variables are the parameters on which the overall cost calculation depends.

$$C_{\text{total}} = \sum_{j \in J} x_j \cdot P_j + \sum_{j \in J} \sum_{k \in K} \sum_{(u,v) \in E} f_{jk}^{(u,v)} \cdot w^{(u,v)},$$

The formulation helps in calculating the overall cost of the big data processing in the distributed data centers. Therefore, reducing the electricity cost has received significant attention.

## VIII. EXPERIMENTAL RESULTS

| Job | proc_time | cost | arr_rate |
|-----|-----------|------|----------|
| 4 | 30 | 3000 | 409.0 |
| 3 | 26 | 2600 | 472.0 |
| 2 | 22 | 2200 | 558.0 |
| 1 | 29 | 2900 | 423.0 |
| 5 | 29 | 2900 | 252.0 |

Table 2: Processing results

The results shows the processing time ,cost of the processing and arrival rate of the jobs .As compared to previous techniques it is found from experimental data that the task placements and data resizing mechanism gives a more cost beneficial results .

## ACKNOWLEDGEMENT

In this paper, we mutually concentrate on the data placement, task assignment, data focus resizing and directing to minimize the general operational cost in huge scale geo-dispersed data places for enormous data applications. We first portray the data preparing procedure and infer the normal finish time in shut structure, in view of which the joint streamlining is formed as a Linearization issue. To handle the high computational many-sided quality of comprehending our Linearization, we Formulated it. Through broad examinations, we demonstrate that our joint-enhancement arrangement has considerable point of interest over the methodology by two-stage separate enhancement. A few intriguing phenomena are likewise seen from the test results.

## REFERENCES

[1] "Data Center Locations," http://www.google.com/about/datacenters/inside/locations/index.html.

[2] R. Raghavendra, P. Ranganathan, V. Talwar, Z. Wang, and X. Zhu,"No "Power" Struggles: Coordinated Multi-level Power Management for the Data Center," in Proceedings of the 13th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS). ACM, 2008, pp. 48–59.

[3] L. Rao, X. Liu, L. Xie, and W. Liu, "Minimizing Electricity Cost: Optimization of Distributed Internet Data Centers in a Multi Electricity-Market Environment," in Proceedings of the 29th International Conference on Computer Communications (INFOCOM). IEEE, 2010, pp .

[4] Z. Liu, M. Lin, A. Wierman, S. H. Low, and L. L. Andrew, "Greening Geographical Load Balancing," in Proceedings of International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS). ACM, 2011, pp. 233–244.

[5] R. Urgaonkar, B. Urgaonkar, M. J. Neely, and A. Sivasubramaniam, "Optimal Power Cost Management Using Stored Energy in Data Centers," in Proceedings of International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS). ACM, 2011, pp. 221–232.

[6] B. L. Hong Xu, Chen Feng, "Temperature Aware Workload Management in Geo-distributed Datacenters," in Proceedings of International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS). ACM, 2013, pp. 33–36.

[7] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," Communications of the ACM, vol. 51, no. 1, pp. 107–113, 2008.

[8] S. A. Yazd, S. Venkatesan, and N. Mittal, "Boosting energy efficiency with mirrored data block replication policy and energy scheduler," SIGOPS Oper. Syst. Rev., vol. 47, no. 2, pp. 33–40, 2013.

[9] I. Marshall and C. Roadknight, "Linking cache performance to user behaviour," Computer Networks and ISDN Systems, vol. 30, no. 223, pp. 2123 – 2130, 1998.

[10] H. Jin, T. Cheocherngngarn, D. Levy, A. Smith, D. Pan, J. Liu, and N. Pissinou, "JointHost-Network Optimization for Energy- Efficient Data Center Networking," in Proceedings of the 27th International Symposium on Parallel Distributed Processing (IPDPS), 2013, pp. 623–634.

[11] A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag, and B. Maggs, "Cutting the Electric Bill for Internet-scale Systems," in Proceedings of the ACM Special Interest Group on Data Communication (SIGCOMM). ACM, 2009, pp. 123–134.

[12] X. Fan, W.-D. Weber, and L. A. Barroso, "Power Provisioning for A Warehouse-sized Computer," in Proceedings of the 34th Annual International Symposium on Computer Architecture (ISCA). ACM, 2007, pp. 13–23.

[13] S. Govindan, A. Sivasubramaniam, and B. Urgaonkar, "Benefits and Limitations of Tapping Into Stored Energy for Datacenters," in Proceedings of the 38th Annual International Symposium on Computer Architecture (ISCA). ACM, 2011, pp. 341–35