# COST MINIMIZATION TECHNIQUES FOR BIG DATA PROCESSING IN DATA CENTERS

**Vasantham Vijay Kumar[1]**

[1]Assistant Professor, Koneru Lakshmaiah Education Foundation, Guntur, India-522302

Mail id: [1] *vijaykumarvasantham@kluniversity.in*

**Abstract:**

In commonly we are using data as single cluster to store in central approaching and its moving with inefficient or infeasible due to lack of different imitations as lack of wide area transmission capacity and the low latency essential of data processing. In big data processing across all the information processing to specific geo-distributes centers working and computing resources. Whatever we are managing distributed data using Map Reduce functions or computations are data will distributed on datacenters not really to solve different types of technical issues. In data centers how to appropriate data surrounded by selection of geo-distributes to minimize the transmission cost, how to resolve the VM (Virtual Machine) managing approach that offers very high performance and low cost this working criteria will select datacenters as outcome of the reduce for big data analytics jobs.

In this paper, these difficulties is tended to by adjusting bandwidth cost, stockpiling cost, processing cost, migration cost, and latency cost, between the two Map Reduce stages across datacenters. We figure this intricate cost streamlining issue for information development, asset provisioning and reducer choice into a joint stochastic number nonlinear improvement issue by limiting the five cost factors all the while.

**Keywords:** Big Data, cost minimization, fuzzy tokens, Mixed-Integer Linear Programming (MILP)

## 1. INTRODUCTION

When the large information blast to adapting and investigate excessive quantities of endless information streams progress. In general the web based communicating stream data, sensor

information stream, log data, stock marketing or trading data and so on., these are applications become a critical needed for some more logic and conditions and tools applications in recent times. In spite of spitefulness of the real that VM positions in geo-dispersed working with server based system and it farms for cost minimization has been commonly contribute to very good effectively applicable to clump adapting models for Map Reduce.

However none of them we can apply to emitting work process below points for realties mention.

1) They oversight to catch the quality of charges in overflow work processing for the short life pattern of all data stream, so they can be effectively to activated to bunch figuring model with the element of establish to all stream information at that point registering them. Sometimes data stream all the information doesn't handling reason of the need to carry out low idleness.

2) In a large data depends on the uncertainty that the cost of VMs and traffic among the servers farms are static and constant while allotting work process into geo-appropriated server farms,

In this paper mainly aim to propose a work process to reduce or minimization of cost with continuous process to handling different data stream of both VMs and traffic.

In Map Reduce is a good convey programming method for preparing hug extension dataset in equivalent. It is also more extraordinary activity in many current applications. Since particular Map Reduce model doesn't improve for system across datacenters. The data to assemble propagate information to a solitary datacenters for integrated approaching is a commonly utilized methodology.

In this process may suspended tight for like a unified aggregate more experience a particular delays the main reason of the heterogeneous and particular data will transfer capacity of user cloud interface. Notice that the data transfer speed of among the data center connectivity is typically carried generally high-transmission capacity lines, when the data will be transforming to numerous data centers for map activity in equal and after that conglomerating the center of the road messages to a solitary data center for reduce activity and utilizing data centers inference can probably diminish the dormancy.

Beside the various types of cost (e.g., acquired by moving the information or rent out VM) and more added can be improved considering the heterogeneity of the linking speed , the dynamism of the data age and the asset cost. In this way, data to spreading circulating information form rotating information form multi-sources into datacenters and ready to preparing them using conveyed map reduce is attention way to manage the extensive volume scattered messages.

As far as now, the most symbolic inquiries to be deal include:

1) How to the position of streamline of extensive scope data sets from contrasting areas onto geo-circulated data center in cloud handling.

2) What are the number of assets, for example, data processing assets ought to be provisioned to ensure compile and accessibility when going limiting the expense.
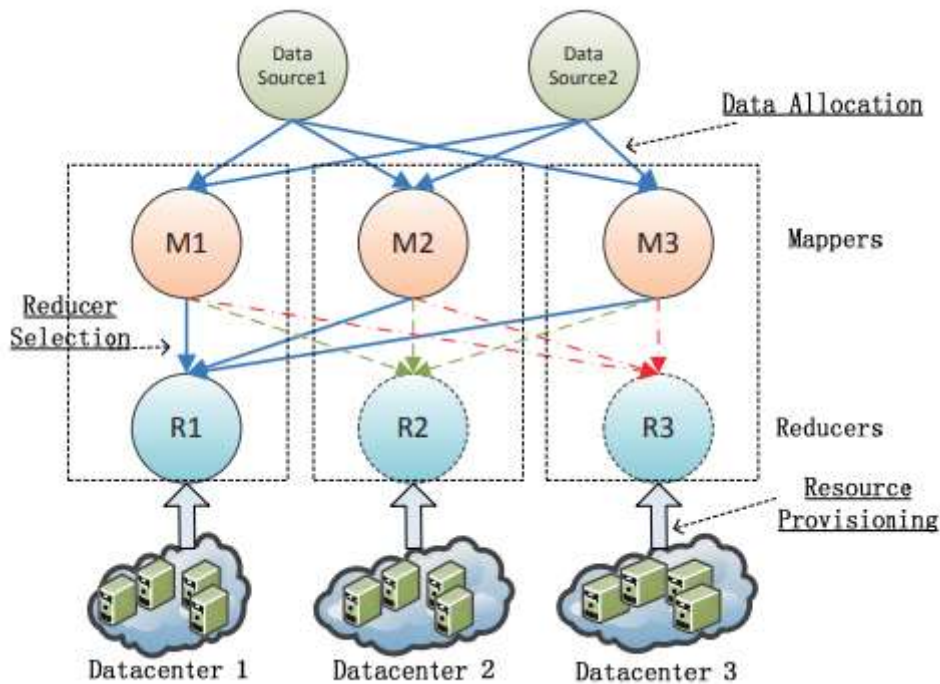
In a dynamic-driven evaluating model a various fluctuation and big data wellsprings of produced information to be merged with cloud asset to give a highly testing problem. The main difference between map stage and the reduce period of map reduce program. It raises the multilevel or multidimensional nature of the information to design and developed. It is provisioning and final data decline choice issues in geo-circulated data centers.

In this paper we discuss about the problem of efficient of a particular objective of superior.

As shown in Fig.1, It is highly accessibility and cost minimization by accommodate five arrange of cost between the three map reduce stages over numerous geo-circulated data centers like as capacity cost, stockpiling cost, registering cost, movement cost, and idleness cost.

We propose  a design can purposely the major draw backs of information developing and introduce a map reducer technique apply on issue to determine within the setting of running over the data centers overall different data centers and VMs of different sorts and randomly cost. We describe not predictable of cost issue as a jointed stochastic total number of non linear improvement streamlining structure by changing the first issue divided into three self determining sub problems that can be settled with some basic arrangements. So we plan an good efficient and advertize online calculation to limit the drawn out of time in the middle value of activity cost.

Fig. 1: Architecture of Big Data Processing with MapReduce Across Datacenters



We consider the explain if Min BDP as far as minimize cost and most discouraged scenario delay. We display that the compute roughly the ideal order within certain limits and secure that the data approach can be finished innards pre-defined delays. We use direct board analysis to assess the proposal of our online calculation with real world data sets. The experiment result displays its competence also its control regarding cost. In this frame work dependent and dynamic time to  actual customer draws near for e.g the blends of data or information distribution systems and the income provisioning of strategies .

## 2. Related works

  The virtual machines already existed we could go to the bottle and that changes the options. The most important all those of virtual machines at once this is where all the virtual machines to be decided to automatically shut down. Every night in the stick the number of it together shutting down the virtual machines, but we cannot shut down this virtual machine that if the

number of different times committed the rigid instances. So when we know the server always to be running it should commit to specify the location and for a long period of a time.

The ability and reservations will find the documentation that explains all of them by reserving. The virtual machines for a long period using it are so very good way to save the money on server virtual machines. A long period of time so the another type the virtual machine the reservation of virtual machines sometimes the little shocked to run the would not want to the way to learn much time. The virtual machines available so that while the stop virtual machines under any process to know job process that you have to so many next step. The maximization of benefits of it the hybrid in the information on the virtual machines.

The lot of information that we assume to the location all the information available by the feature is available in each region. It is very available adjust the between the two regions and needed to hear we want yet the feature is available in each reasons. It is available to adjust between the reasons and needed to hear we want to hear that it's it is a good services. It is a to migrate or migrated for example it is a properly bad ideas we can to multiple servers of just a tool using existing data bases or data process.

**Task Distribution and VM Placement**

The major aim of data distribution and VM arrangement are two different approach to increase the DC execution and usage. Liu et al. [1] propose the Green cloud design Which entitle complete internet checking, live VM movement, and VM situation streamlining Meng et al. [2] develop the system adaptability by buildup the traffic-mindful VM position as per the traffic designs among VMs. Ajiro et al. [3] propose a VM situation conspire called FFD (First-Fit-Decreasing).

The basic principal of cloud server that can be completely effect its gain needed. Zeng et al. [4] address the major issue of VM strategies to limit the assemble related cost inside a DC below the thought of both design and gain requirements. As shown on below Figure2, there are two producers(s1 ad s2) as sources, two consumers(d1 and d2) as destinations and four VMs for four different tasks(t1 to t4). So totally we are use unique in common to real work on VM strategies all permitting that the between VM traffics to consider VM status and stream adjusting.
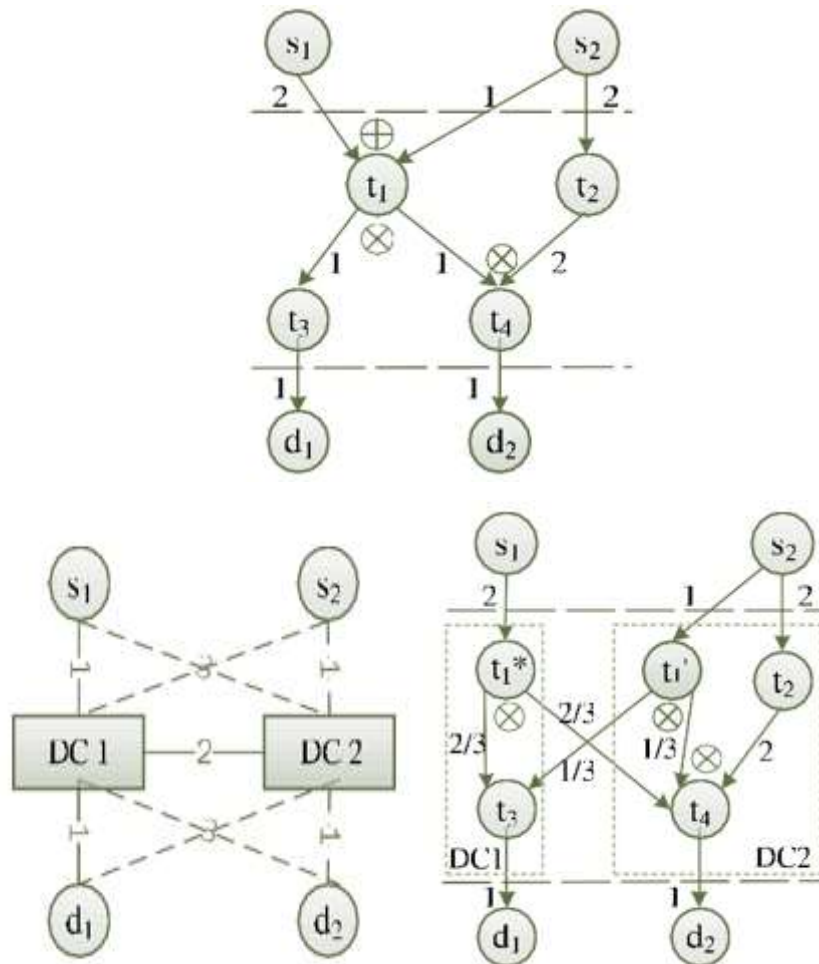
**Fig.2   a) Task Flow                              b) DCs with diverse inter DC network cost**

**c) Flow Graph for Placement**

**Workflow Allocation Algorithm**

From the perspective of target limits, VM placement estimations can be requested into six classes.

In the first place, the works [2], [3] proposed imperativeness gainful booking computations with the objective limit of essentialness consumption minimization, and number of VMs minimization, independently.

.

Second, the works [4], [5] proposed compose traffic headway computations with the objective limit of framework cost-minimization, and framework traffic minimization, separately.

Third, the works [6], [7] proposed reasonable smoothing out figurings with the objective capacity of traditionalist pay development, and operational cost-minimization, independently.

Fourth, the works [8], [9] proposed execution support figurings with the objective limit of openness maximization, and QoS(Quality of Services) enlargement, independently.

Fifth, the works [10], [11] proposed resource utilization maximization figurings with the objective limit of most prominent ordinary use minimization, and resource use support, exclusively.

Finally, the work[12] proficient spoke to a multi-target approach and the target work was planned as a straight mix of various objectives.

Regardless, most computations have been created for bunch data planning, (for example, Map/Reduce model), and the issue of VM portion for stream immense data preparing has not been properly tended to. In this work, we revolve around the cost-minimization of SW designation issue, considering the traits of SW(Software) and the worth heterogeneity of geographically scattered datacenters.

For the data task courses of action, three operator strategies are under the following:

1) Proximity-careful Data Allocation (PDA), in which intensely made data from each data source are continually administered to the topographically nearest datacenter. It produces irrelevant inertia and is proper for the circumstance that inaction delay is before various components.

2) Load-Altering Data Allocation (LADA), in which the data from each data source are continually dispatched to the datacenter with the least Map remaining weight. Obviously, this procedure is good for keeping exceptional main job balanced among datacenters.

3) Minimal Price Data Allocation (MPDA), in which the data from each data source are allotted to the most money related datacenter, to achieve the least cost.

For the VM provisioning approaches, two average strategies are below points:

1) Heuristic VM Provisioning (HVP), in which the VMs required at current time are assessed dependent on the outstanding task at hand at past time. To adapt to the variance of remaining task at hand, additional 50 percent VMs are added to those need at past time to shape an official conclusion.

2) Stable VM Provisioning (SVP), in which the VM include of each kind in each datacenter is set to a fixed worth. For simplicity of correlation, we arrange the fixed an incentive as the normal VM of each kind accomplished by Mini BDP. In this manner, the measure of VMs devoured by SVP is equivalent to that of Mini BDP inside timeframe T .

## 3. Conclusion

In this paper, we investigate the correspondence cost minimization for BDSP ( Big Data Stream Processing)  in geo-scattered DCs by methods for researching the between DC traffic cost tolerable assortments. A MILP (**Mixed-Integer Linear Programming)**  enumerating is worked for this issue, where VM course of action and stream modifying are commonly thought of. To deal with the computational multifaceted nature, we by then propose the "MVP" count subject to our MILP enumerating. Through expansive examinations, we show that our "MVP" computation performs incredibly close to the ideal course of action and inside and out beats the single-VM based BDSP.

## References

[1] W. Andreas, "Data Workflow – A workflow model for continuous data processing", Univ. of Twente, Netherlands, 2010.

[2] M. Tang and S. Pan, "A hybrid genetic algorithm for the ener-gy-efficient virtual machine placement problem in data centers," Neural Processing Letters, vol. 41, no. 2, pp. 211-221, 2015.

[3] M. Sun, W. Gu, X. Zhang, H. Shi, and W. Zhang, "A matrix transformation algorithm for virtual machine placement in cloud," IEEE Proc. 12[th] Int. Conf. TrustCom, pp. 1778-1783, 2013.

[4] R. Wang, J. A. Wickboldt, R. P. Esteves, L. Shi, B. Jennings and L. Z. Granville, "Using empirical estimates of effective band width in network-aware placement of virtual machines in data-centers," IEEE Trans. On Network and Service Management, vol. 13, no. 2, pp. 267-280, 2016.

[5] W. Chen, I. Paik, and Z. Li, "Topology-aware optimal data placement algorithm for network traffic optimization" IEEE Trans. On Computers, vol. 65, no. 8. Pp. 2603-2617, 2016.

[6] W. Shi, and B. Hong, "Towards profitable virtual machine placement in the data center," in 4[th] IEEE Int. Conf. on Utility and Cloud Computing (UCC), pp. 138-145, 2011.

[7] J. J. Prevost, K. Nagothu, B. Kelley, and M. Jamshidi, "Optimal update frequency model for physical machine state change and virtual machine placement in the cloud," in IEEE 8[th] Int. Conf. on System of Systems Engineering (SoSE), pp. 159-164, 2013.

[8] E. Bin, O. Biran, O. Boni, E. Hadad, E. K. Kolodner, Y. Moatti, and D. H. Lorenz, "Guaranteeing high availability goals for virtual machine placement," in IEEE Proc. 31th Int. Conf. on Distributed Computing Systems (ICDCS), pp.700-709, 2011.

[9] K. Sato, M.Samejima, and N. Komada, "Dynamic optimization of virtual machine placement by resource usage prediction," in IEEE 11[th] Int, Conf. on Industrial Informatics, pp. 86-91, 2013.

[10] Y.Guo, A. L.Stolyar, and A. Walid, "Shadow routing based dynamic algorithms for virtual machine placement in a net-work loud," in Proceedings IEEE INFOCOM , pp. 620-628, 2013.

[11] F. Song, D. Huang, H. Zhou, H. Zhang, and I. You, "An optimization based scheme for efficient virtual machine placement," Journal of Parallel Programming, vol. 42, no.5, pp. 853-872, 2014.

[12] S. H. Wang, P. P. W. Huang, C.H.P. Wen, and L.C. Wang, "EQVMP: Energy-efficient and QoS-aware virtual machine placement for software defined datacenter networks," in Proc. 16[th] Int. Conf. on Information Networking, pp. 220-225, 2014.

[13] Maddumala, V.R., Arunkumar, R. (2017). Big data-driven feature extraction and clustering based on statistical methods. Traitement du Signal, Vol. 37, No. 3, pp. 387-394.

[14] Danapaquiame, N. Balaji, V. Gayathri, R. Kodhai, E. Sambasivam, G.(2018). Frequent Item set Using Abundant Data on Hadoop Clusters in Big Data, Vol. 11,No.1, pp. 104-112.

[15] Sree Ram N., Krishna Prasad M.H.M., Satya Prasad K. (2019), 'Repartitioned optimized K-mean centroid based partitioned clustering using mapreduce in analyzing high dimensional big data', International Journal of Engineering and Advanced Technology, 9(1), PP.616-620.

[16] Radhika D., Aruna Kumari D. (2018),'Misusability measure based sanitization of big data for privacy preserving MapReduce programming',International Journal of Electrical and Computer Engineering,8 (6),PP. 4524-4532

[17] Ramakrishna B., Nagabhushana Rao M., Pittala R.B. (2018),'Low cost geo distributed data centers with big data process',Journal of Advanced Research in Dynamical and Control Systems,10 (7),PP. 394-397