



An improved method for allocating cloud resources using MLP in accordance with SLA agreements and user requirements

Dr. P.Madhuri | Chintolla Surekha | Shaik Meer Subhani Ali | V. Nava Kishore

Department of Computer Science and Engineering, Hyderabad Institute of Technology and Management, Hyderabad, Telangana, India.

To Cite this Article

Dr. P.Madhuri, Chintolla Surekha, Shaik Meer Subhani Ali, V. Nava Kishore, An improved method for allocating cloud resources using MLP in accordance with SLA agreements and user requirements International Journal for Modern Trends in Science and Technology, 2024, 10(03), pages. 313-317. <https://doi.org/10.46501/IJMTST1003054>

Article Info

Received: 15 February 2024; Accepted: 02 March 2024; Published: 09 March 2024.

Copyright © Dr. P.Madhuri et al.; This is an open access article distributed under the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

ABSTRACT

Cloud computing is the most popular and important area of study in the world today. The cloud's most active areas are resource provisioning and resource management, which have resulted in a huge variety of solutions to meet these needs. Because of the constraints imposed by users and service level agreements (SLAs), cloud resource allocation is a complicated procedure. In this work, we concentrated on a wide range of problems, ranging from user requests to service level agreements, and suggested a machine learning-based method to address them. Previous methods have only relied on statistical methodologies; machine learning, on the other hand, is an optimum computing methodology that may be used to solve complicated problems efficiently. When compared to the prior methods, the suggested methodology yields the most favourable outcomes.

Keywords: SLA, Cloud computing, user requests, resource allocation, MLP

1. INTRODUCTION

Cloud computing is one of the trendy expression terms in the present world. Practically all Information and Communication Technology (ICT) frameworks are presently moving into cloud based computing model as opposed to conventional computing framework [1]. It is such another computing worldview giving programming, framework, and stage as administrations on-request premise over the Internet. In addition, in contrast to conventional computing framework, cloud

computing requires insignificant administration exertion or specialist co-op communications.

Cloud condition thinking about anticipated execution and SLA (Service Level Agreement) guarantees. Various methodologies, for example, improved burden adjusting calculations, moving the heaps among servers, or changing over the servers into vitality sparing modes (i.e., rest/rest, inert, dynamic/on, and off states) have proposed.

Googles Green Data Centers report [4] recommended three prescribed procedures and five-advance

methodologies for cooling and diminishing vitality use inside data centres. Notwithstanding, resource task in a cloud domain can be practiced in an assortment of ways, for example, proactive administration, responsive administration, etc. Prescient or proactive resource portion is a standout amongst the most dominant and promising methodologies for overseeing cloud resources. It powerfully conjectures and makes the relationship among applications QoS targets, vitality proficiency target capacity, and current equipment resources portion and client remaining task at hand examples.

Thinking about the prescient portion destinations and remaining burden gauging techniques, this paper focuses to build up a cloud model for guaranteeing vitality successful administration situated resource the executives strategy. In the wake of considering different written works, couple of significant issues have been made sense of, i.e.,

(a) portraying the remaining task at hand before forecast; (b) observing atypical resource demands that may damage SLA; (c) versatile expectation instrument with changing the outstanding burdens; and (d) coordinating different resource allotment strategies in a single edge. These issues are incorporated into the proposed model in planning prescient cloud resource distribution model. The remainder of the composition is sorted out as follows. In area II, we present related research foundation with respect to vitality utilization methodologies dependent on outstanding burden expectation. Segment III speaks to the proposed model design dependent on burden forecast. In segment IV, we present how we tested of this work. In area V, we investigate and examine the test results focusing to the objective of the paper. At last, we finish up by recommending future headings in area VI.

2. LITERATURE WORK:

A decade ago, in the early 1980s, researchers began investigating market-driven resource allocation [8][5]. Unlike most other types of market-based resource allocation systems, First Price [3] and First Profit [6] are developed for a set amount of resources and do not use price as a basis for allocation [10]. Our research is concerned with resource allocation for SaaS providers based on user-driven SLA-based profit maximisation. Predictive systems, such as those developed by Reig G.

et al. [11], have made significant contributions to reducing resource use while fulfilling requests and completing them before the deadline. As a result of their prediction system, their scheduling rules might refuse to provide services to requests when the available resource capacity is unable to finish the request before the deadline. Enterprise applications, as opposed to computing and scientific applications, are the focus of our research. Fu Y. et al [21] introduced a dynamic scheduling technique (Squeeze) for streaming distributed resources that was based on service level agreements (SLAs). Yarmolenko V. et al [22] also conducted an assessment of several SLA-based scheduling algorithms on parallel computing resources, using resource (number of CPU nodes) usage and revenue as the evaluation criteria for their findings. Our research, on the other hand, is focused on the scheduling of corporate applications running on virtual machines in Cloud computing settings. (In our work, the smallest unit of resources is the number of virtual machines.) A number of QoS elements on the resource provider's side, such as price and proposed load, were taken into consideration by Popovic et al. [6], but the user side was not taken into consideration]. However, our suggested work varies in terms of QoS parameters from both the customer's and the SaaS provider's perspectives, and it focuses on user-driven situations rather than business-driven ones. In their research, Lee et al. [2] looked at profit-driven service request scheduling for work flow. Our study, on the other hand, a) focuses on SLA-driven QoS parameters on both the user and provider sides, and b) addresses the difficulty of dynamically changing customer requirements in order to increase profit and reputation.

Using genetic algorithms in virtualized settings, Song et al. [18] discussed resource allocation algorithms for enterprise applications in the context of resource allocation methods for corporate applications. Genetic algorithms, on the other hand, need a significant amount of time to execute. In cloud computing settings, where consumers expect to be supplied instantly, the extended execution time raises the likelihood of SLA violations.

3. PROPOSED WORK:

Cloud resource allocation was usually provisioned as per the agreement called SLA of demands and services providers allocates resources by considering the QoS.

This paper provides an architecture for cloud resource allocation by considering the user requests, SLA agreement and availability of services. here we use machine learning mechanism for predicting the load distribution, resource allocation performed by controller. Remaining burden classifier is discretionary and a piece of burden indicator. It orders the heap into high, moderate, or low classes dependent on the anticipated resource use (e.g., CPU usage) for a particular timeframe with some predefined rules. For instance, at t1 and t2 time interim, in the event that the heap indicator appears (80-90)% CPU utilization, at that point it tends to be high CPU load. The resource chief chooses the amount of resources necessity for the approaching burdens. The choice either arrangement or discharge relies upon the outcomes from burden indicator, SLA no tifier, and current resource data. Moreover, resource director coordinates different improvement strategies, for example, green booking calculation [14], or vitality mindful allotment calculation [15] for fine-grained flexible resource the executives. At long last, the heap follow screen consistently feds history outstanding burden follows into expectation model after the preparing of approaching solicitations by cloud servers. At whatever point new demands are overhauled appropriately, the heap follow screen accumulates and passes follows to the indicator.

Resource pool: it contains all kinds of resources with multiple instances.

Monitor: monitor is a machine which monitors all the available resources and give the information to the controller regarding available as well as allocated resources.

Allocator: Responsibility of allocator is to allocate desired resources for users based on the inputs from the controller.

Controller: controller is the key part in our mechanism, he can take care about monitor, allocator, resource pool, SLA, SLA DB and user.

SLA: It is the agreement digital copy which stored in SLA DB.

SLA DB: it stores all users SLA agreements in a prescribed format

Users: who seeks there sources from cloud?

Algorithm for resource allocation based on SLA and user requests:

```

Controller()
Capacity Planning and Auto-scaling()
{
Input: Resource Dem and Utilization of VMS presently
and current resources
Output: Decision on Capability Planning and Auto-
scaling
Notations: WebVM-i: VM running Transactional (Web)
Applications;
CurResDemand(WebVM-i): Current Resource
Demand; CurAllocResWebVM-i: Current Allocated
Capacity; Reserved Res(WebVM-i):
Reserved VMs Capacity Specified in SLA; HpcVM-i:
VM running HPC Application 1: for
Each Web VM-i do
2: Calculate the present asset request Cur
Res Demand (WebVM-i)
3: on the off chance that
CurResDemand(WebVM-i)
<CurAllocResWebVM-iatthatpoint
4: ReducetheassetlimitofWebVM-itocoordinate the
interest
5: else
6: intheeventthatCurResDemand(WebVM-i)≤
ReservedRes(WebVM-i) at that point
7: IncreasetheassetlimitofWebVM-itocoordinate the
interest
8: Reduce correspondingly the asset limit
distributed to HPC application (HpcVM-i
)onasimilarserver
9: else
10: ontheoffchancethatSLAcontainsAuto- scaling
Option, at that point
11: InitiatenewVMsandoffloadtheapplication request to
new VMs
12: end
13: end

```

```

14:end
15:end
16:forEachBatchJobHpcVM-ido
17:onthef chance that slack assets accessible on the
server where HPC VM is running, at that point 18:
Allocate the slack assets
19:endif
20: Recomputed the evaluated completion time of the
activity
21:RescheduletheBatchJobVMifmissingthe due date.
22:endfor

```

Algorithm(forwardpass):

Require: pattern~x, MLP, enumeration of all neurons in topological order

Ensure: calculate output of MLP 1: for all input neurons i do

2: set $a_i \leftarrow x_i$

3: endfor

4: for all hidden and output neurons i in topological order do

5. $set\ net_i \leftarrow w_{i0} + \sum_{j \in Pred(i)} w_{ij} a_j$

6. $set\ a_i \leftarrow f_{top}(net_i)$

7: endfor

8: for all output neurons i do

9: assemble i in output vector ~y 10: end for

11: return ~y

Algorithm2:(forwardpass):

Proposed machine learning based mechanism used for to get the allocation of resources based on SLA and need of users. This should be input for the controller. And controller resource allocation based machine learning model. will take care about on the input of the Here it includes two phases in initial phase training of MLP with the incoming data of resource requests and SLA phases in initial phase training of MLP with the incoming data of resource requests and SLA agreements. After training testing of the mechanism will be happened

4. RESULTS AND DISCUSSIONS

Order Accuracy is the thing that we normally mean, when we utilize the term accuracy. It is the proportion of number of right expectations to the complete number of info tests. computation time is proportional to the number of rule applications

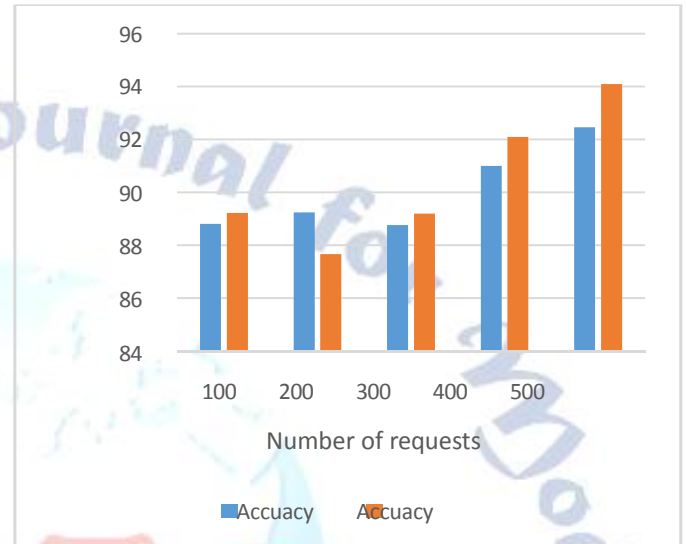


Fig-2: Accuracy of allocation

Error rate

Error rates refer to the frequency of errors occurred, defined as “the ratio of total number of data units in error to the total number of data units transmitted.

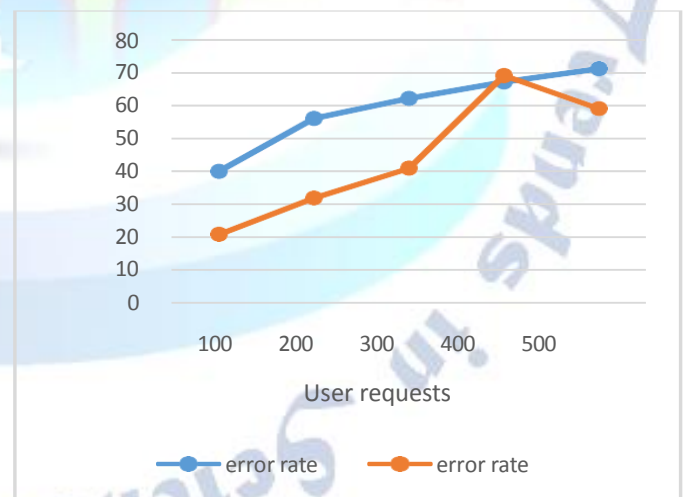


Fig-3: Error rate in allocation of resources with SLA

Computation time

Computation time (also called "running time") is the length of time required to perform a computational process. Representation a computation as a sequence of rule applications, the computation time is proportional to the number of rule applications



Fig-4:Computation time in allocation of resources with SLA

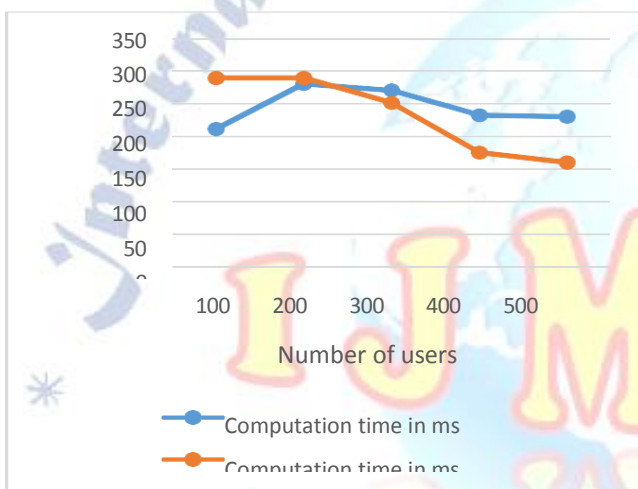


Fig-5:Response time in allocation of resources with SLA

5. CONCLUSION:

Customers have access to three main kinds of on-demand services in cloud computing environments: software as a service, infrastructure as a service, and platform as a service. Customers' requirements for SaaS providers were the subject of this research, which aimed to reduce costs by addressing dynamic needs. In order to do this, we used mapping and scheduling algorithms to cope with customer-side dynamic needs and resource level heterogeneity, as discussed in the introduction. Three algorithms were then constructed that took into account several quality of service (QoS) metrics, such as arrival rate, service commencement time, and penalty rate from both consumers' and SaaS providers' perspectives. The ProfminVMminAvaiSpace method, when compared to the other offered algorithms, optimised cost reductions better on average, according to the simulation findings. In the future, we will look at

ways to increase the profitability of the algorithms in terms of total profit, and we will also look at the SLA negotiation process in Cloud computing environments in order to improve customer satisfaction.. Additionally, we'd want to include new service offerings and pricing tactics like spot pricing to help service providers make a bigger profit. To further enhance our algorithms' time complexity, we're looking at knowledge-based scheduling for optimising a SaaS service provider's profit. In addition, we'll examine the penalty limit in light of system failures.

Conflict of interest statement

Authors declare that they do not have any conflict of interest.

REFERENCES

- [1] S. Yeo, and R. Buyya, "Service level agreement based allocation of clusterresources: Handling penalty to enhance utility". In Proceedings of the 7th IEEE International Conference on Cluster Computing (Cluster2005), Boston, MA, USA.
- [2] Y.C.Lee,C. Wang,A.Y. Zomaya and B.B. Zhou,"Profit-drivenServiceRequestScheduling inClouds".InProceedingsoftheInternational Symposium on Cluster and Grid Computing,(CCGrid2010), Melbourne,Australia.
- [3] O. F. Rana, M. Warnier, T. B. Quillinan,F. Brazier,andD.Cojocarasu,"ManagingViolations inServicelevelagreements".Inproceedings of the5thInternationalWorkshoponGrid Economicsand Business Models (GenCon 2008), Gran Canaris, Spain.
- [4] D.E. Irwin,and L.E. Grit, and J.S. Chase, "Balancing Risk andReward in a Market-based Task Service". In Proceedings of the 13th International Symposium on High PerformanceDistributedComputing (HPDC 2004), Honolulu, HI, USA.
- [5] Y. Yemini, "Selfish optimization in computer networksprocessing".In Proceeding of the 20th IEEE Conference on Decision and Control including the Symposium on Adaptive Processes, San Diego, USA.
- [6] I.Popovici,andJ.Wiles,"Proitableservices in an uncertain world". In Proceeding of the18th Conference on Supercomputing (SC 2005), Seattle, WA.
- [7] R.Buyya,C.S.Yeo,S.Venugopal,J. Broberg,andI.Brandic,"CloudComputing and Emerging IT Platforms: Vision, Hype, and Reality for Delivering Computing as the5thUtility,FutureGenerationComputerSystems",25(6),(pp.599-616),ElsevierScience, Amsterdam, The Netherlands.
- [8] D. Parkhill,"The challenge of the computer utility", 1966, Addison-Wesley Educational Publishers Inc., USA.
- [9] M.A.Vouk,"CloudComputing- Issues, Research and Implementation". In Proceedings of 30th International Conference on Information echnology Interfaces (ITI 2008), Dubrovnik, Croatia.