

International Journal of Computer Science and Mobile Computing



A Monthly Journal of Computer Science and Information Technology

ISSN 2320-088X

IJCSMC, Vol. 3, Issue. 4, April 2014, pg.928 – 933

REVIEW ARTICLE

A Review on Virtual Machine Scheduling in Cloud Computing

Archana Pawar¹ (M.Tech Scholar), Prof. Deepak Kapgate²

^{1,2}Department of C.S.E., GHRAET, Nagpur, Nagpur University

¹amru.vpawar@gmail.com, ²deepakkapgate32@gmail.com

Abstract- Cloud computing is a model for enabling ubiquitous, convenient, on demand network access to a shared pool of configurable computing resources like network, server, storage, applications and services. Since these resources can be rapidly provisioned and released with minimal management effort or service provider interaction, we can see nowadays there is a rapid increase in the use of cloud computing. Thus, it has become an important issue to balance the cloud and its resources so as to provide better performance and services to the end users of the cloud and at the same time majority of users being served by application deployments in cloud provider's environment. In cloud computing, load balancing means balancing three important stages through which a request is processed. This includes Data center selection, Virtual Machine Scheduling and Task Scheduling at selected data center. In this paper, we mainly focus on various techniques and algorithms available for Virtual Machine management. It also gives a clear insight of their characteristics to resolve accumulated load in an efficient Virtual Machine Management.

Keywords- Cloud Computing, Data center, Virtual Machine, VM Scheduling Algorithms, load balancing

I. INTRODUCTION

Cloud computing is an on demand service in which shared resources, information, software and other devices are provided according to the clients requirement at specific time. It can be viewed as solution where data storage and any processing take place without the user being able to pinpoint the specific computer carrying. Cloud computing insures access to virtualized it resources that data center are presented and are shared by others. It is common to divide cloud computing into three categories:

A. Infrastructure as a service (IaaS)

It provides flexible ways to create use and manage virtual machines. In IaaS model, computing resources such as storage, network, and computation resources are provisioned as services. Consumers are able to deploy and run arbitrary software, which can include operating systems and applications. Consumers do not manage or control the

underlying cloud infrastructure but have to control its own virtual infrastructure typically constructed by virtual machines hosted by the IaaS vendor. This thesis work mainly focuses on this model, although it may be generalized to also apply to the other models.

B. Platform as a service (PaaS)

PaaS provides all the resources that are required for building applications and services completely from the Internet, without downloading or installing software. It focuses on providing the higher level capabilities more than just virtual machines required to support applications. In the PaaS model, cloud providers deliver a computing platform and/or Solution stacks typically including operating system, programming language execution environment, database, and web server [5]. Application developers can develop and run their software on a cloud platform without having to manage or control the underlying hardware and software layers, including network, servers, operating systems, or storage, but maintains the control over the deployed applications and possibly configuration settings for the application-hosting environment.

C. Software as a service (SaaS)

In SaaS, the user uses different software applications from different servers through the Internet. It is the application that provides business value for users. In the SaaS model, software applications are delivered as services that execute on infrastructure managed by the SaaS vendor. The provider does all the upgrades and patching while keeping the infrastructure running. Consumers are enabled to access services over various clients such as web browsers and programming interfaces, and are typically charged on a subscription basis [6]. The implementation and the underlying cloud infrastructure where it is hosted is transparent to consumers. Eg, Customer resource management (CRM), Video conferencing, IT service management, Accounting, Web content management

Deployment Models in Cloud Computing:

The cloud computing deployment model describes where the software runs and includes the following options: Based on the classification of cloud services into SaaS, PaaS, and IaaS, two main stakeholders in a cloud provisioning scenario can be identified, i.e., the Infrastructure Provider (IP) who offers infrastructure resources such as Virtual Machines, networks, storage, etc. which can be used by Service Providers (SPs) to deliver end-user services such as SaaS to their consumers, these services potentially being developed using PaaS tools. As identified in [7], four main types of cloud scenarios can be listed as follows.

- 1) *Private cloud*: Private cloud is set of standardized computing resources that are dedicated to an organization, usually on-premises in the organization data centers. In private cloud-based service, data and processes are managed within the organization without the restrictions of network bandwidth, security exposure and legal requirements. It works with the current capital investment and drives the new function as a service.
- 2) *Cloud Bursting*: Private clouds may offload capacity to other IPs under periods of high workload, or for other reasons, e.g., planned maintenance of the internal servers.
- 3) *Federated Cloud*: Federated Cloud are cloud collaborated on a basis of load sharing agreements enabling them to offload capacity to each other's in a manner similar to how electricity providers exchange capacity.
- 4) *Multiple clouds*: In multi-cloud scenarios, the SP is responsible for handling the additional complexity of coordinating the service across multiple external IPs, i.e. planning, initiating and monitoring the execution of services.

Parameters of interest for cloud services Provider

- 1) *Resources utilization details*: Just like any other performance monitoring utilization parameter of physical server infrastructure is an important factor in cloud monitoring, as these service make up the cloud.
- 2) *Infrastructure response time (IRT)*: IRT gives the clear picture of the overall performance of the cloud as it checks the time taken for each transaction to complete.
- 3) *Virtualization metrics*: Similar to the physical machine, we need to collect the resource utilization data from the virtual machines. This provides the picture of how much of the virtual machine is being utilized and this data helps in the resources utilization by application and divided on the scale requirements.
- 4) *Transaction matrices*: It can be considered as derivative from IRT . Metrics like success percentage of transaction counts of transaction etc. for an application would give a clear picture of the performance of an application in cloud particular instant.

Cloud computing enjoys the many attractive attributes of virtualization technology, such as consolidation, isolation, migration and suspend/resume support. A virtual machine (VM) is a software implementation of a computing environment in which an operating system (OS) or program can be installed and run. Important parameters related to virtual machines are Number of virtual machines used by applications, Time taken to create a new VM, Time taken to move an application from one VM to another, Time taken to allocate additional resources to VM. Virtualization means “something which isn’t real”, but gives all the facilities of a real. Virtualization is the creation of a virtual version of something such as an operating system, a server, a storage device or network resources.

Scheduling the basic processing units on a computing environment has always been an important issue [1]. Like any other processing unit, VMs need to be scheduled on the cloud in order to maximize utilization, Do the job faster, Consume less energy, Easy resource reservation (allocation). VM’s elasticity in cloud computing, elasticity is defined as the degree to which a system is able to work loud change by provisioning and de-provisioning resources in an automatic manner such that at each point in time the available resources match the current demand as closely as possible.

The number of cloud users has been growing exponentially and apparently scheduling of virtual machines. In the cloud becomes an important issue to analyze. In cloud computing, a user may require a set of virtual machine co-operating with each other to accomplish one task. In the past the inter relationship among task are not considered. Scheduling is the method by which virtual machine flows are given access to system resources.

II. VM SCHEDULING ALGORITHM

The goal of scheduling algorithms in distributed systems is spreading the load on processors and maximizing their utilization while minimizing the total task execution time Job scheduling, one of the most famous optimization problems, plays a key role to improve flexible and reliable systems. The main purpose is to schedule jobs to the adaptable resources in accordance with adaptable time, which involves finding out a proper sequence in which jobs can be executed under transaction logic constraints [2].

There are main two categories of scheduling algorithm. First is static scheduling algorithm and another is dynamic scheduling algorithm. Both have their own advantage and limitation. Dynamic scheduling algorithms have higher performance than static algorithm but have a lot of overhead compare to it. Decisions on load balancing are based on current state of the systems. It is better than static approach. Static algorithms are mostly suitable for homogeneous and stable environments and can produce very good results in these environments. It doesn’t depend on the current state of the system. Prior knowledge of the system is needed. However, they are usually not flexible and cannot match the dynamic changes to the attributes during the execution time. Dynamic algorithms are more flexible and take into consideration different types of attributes in the system both prior to and during run-time. These algorithms can adapt to changes and provide better results in heterogeneous and dynamic environments. However, as the distribution attributes become more complex and dynamic. As a result some of these algorithms could become inefficient and cause more overhead than necessary resulting in an overall degradation of the services performance.

A) Gang scheduling Algorithm

C. Reddy [7] explain use of gang scheduling algorithm in cloud computing responsible for selection of best suitable resources for task execution, by taking some static and dynamic parameters and restrictions of VM into the considerations. Gang scheduling is a scheduling algorithm for parallel system that scheduled related VM to run simultaneously on different machines. Gang Scheduling is an efficient job scheduling algorithm for time sharing, already applied in parallel and distributed systems. Gang scheduling can be effectively applied in a Cloud Computing environment both performance-wise and cost-wise. Gang scheduling is a special case of job scheduling that allows the scheduling of such virtual Machines. Gang scheduling is a special case of scheduling parallel jobs in which tasks of jobs need to communicate very frequently. Gang scheduling involves high overhead since network status must be saved and then be restored when switching between tasks.

Moschakis et. al. [8] gives improved version of gang scheduling and performance and Cost evaluation of Gang Scheduling. Usually, the scheduling methods implemented in the scheduler aim for better response times and lower slowdowns, by minimizing unnecessary delays. The scheduler must also tend to the cost of the lease time of VMs aiming for a better cost-to-performance ratio [8]. The study takes into consideration both performance and cost while integrating mechanisms for job migration and handling of job starvation. The number of Virtual Machines (VMs) available at any moment is dynamic and scales according to the demands of the jobs being serviced. Results highlight that this scheduling strategy can be effectively deployed on Clouds. The number of VM’s available at any moment is dynamically scales according to the demand of the job being served. For this purpose they applied “shortest Queue First (SQF) algorithm which dispatches the tasks to VMs with the shortest queue.

The research on gang scheduling has shown the potential of time sharing in improving throughput [17]. Migrating VM to a new data center is generally expensive, and also it does not eliminate starvation. Gang scheduling can be effectively applied in cloud computing environment both performance wise and cost-wise.

Response Time R_j of a job j is the time interval between the arrival and the departure of the job. Its average is defined as [8]:

$$\text{Response Time } R_j = \sum_{J=1}^n (R_j / n)$$

Where n is the total number of jobs.

B) Round Robin algorithm

In the round robin scheduling, Virtual Machines are dispatched to physical hardware in a FIFO manner but are given a limited amount of CPU time [6] called a time-slice or a quantum. If a process does not complete before its time quantum, the Virtual Machine execution is preempted and given to the next Virtual machine waiting in a queue. The preempted process is then placed at the back of the ready list. A time quantum is generally from 100-1000 milliseconds. So, the RR algorithm will allow the first VM in the queue to run until it expires its quantum (i.e. runs for as long as the time quantum), then run the next VM in the queue for the duration of the same time quantum. The RR algorithm is naturally pre-emptive. RR algorithm is one of the best scheduling algorithms that developed by many researchers

RR is proportionally fair algorithm, or maximum throughput scheduling (throughput). The main advantage of this algorithm is that it utilizes all the resources in a balanced order (resource utilization). The scheduler starts with a node and moves on to the next node, after a VM is assigned to that node. This is repeated until all the nodes have been allocated at least one VM and then the scheduler returns to the first node again. Hence, in this case, the scheduler does not wait for the exhaustion of the resources of a node before moving on to the next (Fault tolerant) [6].

C) Genetic algorithm

Genetic algorithm is for scheduling sets of independent VM's, the objective of genetic algorithm is to minimize the make span. Initially in GA many individual solutions are (usually) randomly generated to form an initial population. The population size depends on the nature of the problem i.e. type and no of VM's to be run effectively on system. During each successive generation, a proportion of the existing population is selected to breed a new generation. Individual solutions are selected through a fitness-based process, where fitter solutions (VM's schedule likely to give effective response time) are typically more likely to be selected. The next step is to generate a second generation population of solutions from those selected through genetic operators: crossover and mutation. This generational process is repeated until a termination condition has been reached i.e. a solution is found that satisfies minimum response time criteria.

GA will increase the cost of time, space, throughput and improve the quality of service of the entire .The goal of GA is to reduce the scheduled time of VM. Genetic algorithm provides both improved response time to VM via parallel execution. A state of the system and through genetic algorithm the migration cost becomes a problem.

D) Content-Based Virtual Machine Scheduling Algorithm

The content based VM scheduling algorithms were designed with the goal of lowering the amount of data transferred between racks in the data center when virtual machines disk image are being copied to the host node [5]. The algorithm returns the selected node and the VM on that node with the highest similar content. When deploying a VM, we search for potential hosts that have VMs that are similar in content to the VM being scheduled. Then, we select the host that has the VM with the highest number of disk blocks that are identical to ones in the VM being scheduled.

Once we have chosen that host node, we calculate the difference between the new VM and the VMs residing at the host; then, we transfer only the difference to the destination host. Finally, at the destination host, we can reconstruct the new VM from the difference that was transferred and the contents of local VMs. Content based VM scheduling algorithm that can significantly reduced the network traffic associated with transfer of VM's from storage racks to host racks in cloud data center.

E) Adaptive Algorithm

K. Kumar [6] proposed adaptive algorithm which uses dynamic priority for nodes based on which the VM's are scheduled and assigned. Depending upon the priority values, the VM's to the nodes are scheduled, which varies dynamically based on their load factor. Priority of a node is assigned depending upon its capacity and the load factor. This algorithm strikes the right balance between performance and power efficiency as and when the virtual

machines are assigned to the nodes, recalculation of their priorities takes place. The dynamic Priority concept leads to better utilization of the resources. Adaptive algorithm is an efficient algorithm for finding expected response time of each Virtual machine. It improve the throughput, achieves high bandwidth utilization and outage probability of the system.

F) *Priority scheduling algorithm*

The basic idea is straightforward; each Virtual Machine is assigned a priority, and priority is allowed to run. Equal-Priority instances are scheduled in FCFS order. Priorities are assigned based on the characteristics of VM's such as amount of workload, predicted execution time, user assigned priority. Internally defined priorities use some measurable quantities or qualities to compute priority of a VM. Priority once assigned to VM can be changed dynamically by using concept of aging i.e. here the priority of VM keeps increasing based on the total amount of time VM remains in ready queue waiting for execution. If priority of VM increased greater than the VM executing on physical hardware the executing VM preempts with VM having higher priority. Preemption of VM from physical hardware is also done when a VM is created or migrated to system having higher priority than VM executing on hardware.

Vignesh V et. al. [13] in their paper proposed improved priority scheduling algorithm using SJF policy. The shortest-Job-First (SJF) policy is used as a special case of general priority scheduling algorithm. An SJF algorithm is simply a priority algorithm where the priority is the inverse of the next CPU burst. That is, the longer the CPU burst, the lower the priority and vice versa [13].

It has High processor utilization, High throughput, Minimize turnaround time. It can be change its priority based on its age or execution history.

G) *Efficient Resource Utilization Algorithm*

R. Nivethitha [9] explains use of the massive pool of resources in terms of pay-as-you use policy. On demand the resources are delivers by the cloud through the use of network resources under different conditions. Based on their usage the effective utilization of resources the users will be charged. His named his proposed algorithm as "Effective Resource Utilization Algorithm (ERUA)" is based on 3-tier cloud architecture (Consumer, Service Provider and the Resource Provider) which benefits both the user (QoS) and the service provider (Cost) through effective schedule reallocation based on utilization ratio leading to better resource utilization.

Performance analysis made with the existing scheduling techniques shows that efficient resource utilization algorithm gives out a more optimized schedule and enhances the efficiency rate [9]. The service provider hires resources from the resource provider and creates Virtual Machine (VM) instances dynamically to serve consumers.

H) *Renewable Energy Source provisioned algorithm*

D. Hatzopoulos et. al. [11] in their paper explores the problem of virtual machine (VM) allocation in a network of cloud server facilities which are deployed in different geographical areas [11]. He addresses the problem of energy-efficient allocation in the system. The objective is to reduce the total cost of power consumption for the operator. Each request for a task to be executed in the cloud is associated with a VM request with certain resource requirements and a deadline by which it needs to be completed. The cloud provider has to create a VM with the resource requirements of the request and to execute the VM before the deadline. He proposes an online algorithm with given look-ahead horizon, in which the grid power price and pattern of output power of the RES are known a priori.

Renewable Energy Source provisioned algorithm applications leads to significant reduction of waste in bandwidth resource. This approach and design the metering method to observe the resource states. It determines its maximum request throughput.

III. CONCLUSION AND FUTURE WORK

In this paper, we have studied and surveyed various algorithms for efficient load balancing and managing virtual machine. Each algorithm has their own pros and cons but at the same time gives us clear scenario where they can be most appropriately suited. It also gives insight of effectiveness and various performance characteristics of these algorithms. These algorithms tell us what kind of parameters we should take into consideration while selecting VM. As a future work, we are trying to extend the scope of VM to be applicable for multimedia applications which is emerging technology and lots of users are fascinating towards it. By providing efficient VMM for cloud computing environment will give effective way to manage Multimedia application services and can be efficiently extend to many users.

REFERENCES

- [1] Hadi Salimi , “Advantages, Challenges and Optimizations of Virtual Machine Scheduling in Cloud Computing Environments” in International Journal of Computer Theory and Engineering Vol. 4, No. 2, April 2012.
- [2] Pinal Salot , “A Survey Of Various Scheduling Algorithm In Cloud Computing Environment” in M.E, Computer Engineering, Alpha College of Engineering, Gujarat, India , Volume: 2 Issue: 2.
- [3] MR.NISHANT, “Pre-Emptable Shortest Job Next Scheduling In Private Cloud Computing” in journal of information, knowledge and research computer engineering, NOV 12 TO OCT 13 | VOLUME – 02, ISSUE – 02.
- [4] TARUN GOYAL, “Host Scheduling Algorithm Using Genetic Algorithm In Cloud Computing Environment”, International Journal of Research in Engineering & Technology (IJRET) Vol. 1, Issue 1, June 2013.
- [5] Sobir Bazarbayev, “Content-Based Scheduling of Virtual Machines (VMs) in the Cloud” in University of Illinois at Urbana-Champaign, AT&T Labs Research.
- [6] Kiran Kumar et. al., “An Adaptive Algorithm For Dynamic Priority Based Virtual Machine Scheduling In Cloud” in IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 6, No 2, November 2012.
- [7] Dr. Chenna Reddy , “An Efficient Profit-based Job Scheduling Strategy for Service Providers in Cloud Computing Systems” in International Journal of Application or Innovation in Engineering & Management (IJAIEM) , Volume 2, Issue 1, January 2013.
- [8] I. Moschakis, H. Karatza, “Performance and Cost evaluation of Gang Scheduling in a Cloud Computing System with Job Migrations and Starvation Handling” in Department of Informatics Aristotle University of Thessaloniki, Greece , IEE 2013.
- [9] Ramkumar N, Nivethitha , “Efficient Resource Utilization Algorithm (ERUA) for Service Request Scheduling in Cloud” in International Journal of Engineering and Technology (IJET) , Vol 5 No 2 Apr-May 2013.
- [10] Dimitris Hatzopoulos, “Dynamic Virtual Machine Allocation in Cloud Server Facility Systems with Renewable Energy Sources” at IEEE International Conference on Communications (ICC) 2013, Budapest, Hungary.
- [11] Vignesh V, Sendhil Kumar KS, Jaisankar N , “ Resource management and scheduling in cloud environment” in International Journal of Scientific and Research Publications, Volume 3, Issue 6, June 2013.
- [12] Jeongseob Ahn, Changdae Kim, Jaeung Han, “ Dynamic Virtual Machine Scheduling in Clouds for Architectural Shared Resources”.
- [13] Tommaso Cucinotta, “ Providing Performance Guarantees to Virtual Machines using Real-Time Scheduling” in Tommaso Cucinotta, Dhaval Giani, Dario Faggioli, and Fabio Checconi Scuola Superiore Sant’Anna, Pisa, Italy.
- [14] Junliang Chen, Bing Bing Zhou, “Throughput Enhancement through Selective Time Sharing and Dynamic Grouping” in 2013 IEEE 27th International Symposium on Parallel and Distributed Processing.
- [15] Manoranjan Dash , “ Cost Effective Selection of Data Center in Cloud Environment” in ISSN, Volume-2, Issue-1, 20131.
- [16] Abirami S.P., Shalini Ramanathan (2012), “Linear Scheduling Strategy for Resource allocation in Cloud Environment”, International Journal on Cloud Computing and Architecture, vol.2, No.1, February.
- [17] Nilabja Roy, Abhishek Dubey and Aniruddha Gokhale, “Efficient Autoscaling in the Cloud using Predictive Models for Workload Forecasting”.
- [18] Soramichi Akiyama, Takahiro Hirofuchi, Ryousei Takano, Shinichi Honiden (2012), “MiyakoDori: A Memory Reusing Mechanism for Dynamic VM Consolidation”, Fifth International Conference on Cloud Computing, IEEE 2012.
- [19] V. Venkatesa Kumar et. al , “Job Scheduling Using Fuzzy Neural Network Algorithm in Cloud Environment”, International Journal of Man Machine Interface, Vol. 2, No. 1, March 2012.