

# Dynamic Resource Allocation using Virtual Machines for Cloud Computing Environments

Roneld Khunjan  
M. Tech, CSE student  
T. John Institute of Technology  
Bangalore, India

B. E. Narasimhhaya  
Assistant Professor, Dept. of CSE  
T. John Institute of Technology  
Bangalore, India

**Abstract** – With the recent advancement in cloud computing technology, cloud computing allows users to upgrade and downgrade their resource usage based on their needs. Most of these benefits are achieved from resource multiplexing through virtualization technology in the cloud model. Using the virtualization technology the data center resources can be dynamically allocated based on application demands. The concept of "green computing" and "skewness" is introduced to optimize the number of servers in use and to measure the unevenness in the multi-dimensional resource utilization of a server respectively. By minimizing skewness, the different types of workloads can be combined effectively and the overall utilization of server resources can be improved.

**Keywords** – Cloud Computing, Resource Management, Virtualization, Green Computing.

## I. INTRODUCTION

The flexibility and the lack of capital investment offered by cloud computing is appealing to many businesses. A lot of discussions has been made on the benefits and investments of the cloud model and on how to migrate legacy applications onto the cloud environment. Here a different problem is studied: how can a cloud service provider best multiplex its virtual resources onto the physical hardware? This is important because most of the important benefits in the cloud model come from such multiplexing. It is known from the observation that servers in many existing data centers are not effectively utilized due to over-provisioning of the data center resource usage for the peak demand. The cloud model to be presented is expected to offer automatic scaling up and down of resources in response to load variation to make such practice unnecessary. Besides decreasing the hardware cost, it also reduces the consumption of electricity which contributes to a significant portion of the operational expenses in large data centers.

The mechanism to map virtual machines (VMs) to physical resources is provided by Virtual Machine Monitors (VMMs) such as Xen. And this mapping is made unknown to the cloud users. For example the users with the Amazon EC2 service do not know where their VM instances run. Only to the cloud service provider makes sure that the underlying physical machines (PMs) have sufficient resources to supply their needs. The mapping between VMs and PMs, while applications are running, is made possible by the VM live migration technology [2].

However, a problem still exists as how to decide the mapping effectively so that the resource demands of VMs are met while the number of PMs used is minimized. This becomes a challenging point when the resource needs of VMs are heterogeneous due to the different types of applications they run and update with time as the workloads increase and decrease with time. The capacity of PMs can also be heterogeneous because multiple generations of hardware co-exist in a data center.

The algorithm used focuses on achieving the following two goals:

- Overload avoidance: The existing capacity of PM should have sufficient resources to satisfy the resource needs of all VMs running on it. Else, the PM gets overloaded and may decrease the performance of its VMs.
- Green computing: the number of PMs used should be minimized as long as they can still satisfy the needs of all VMs. Idle PMs can be shut down to save power.

There is an inherent trade-off between the two goals in the face of changing resource needs of VMs. For overload avoidance, the utilization of PMs should be kept low to reduce the possibility of overload in case the resource needs of VMs increase later. For green computing, the utilization of PMs should be kept reasonably high to make efficient use of their energy.

In this paper, the design and implementation of an automated resource management system that achieves a good balance between the two goals is presented and the following contributions is made:

- A resource allocation system is developed that can avoid overload in the system effectively while minimizing the number of servers used.
- The concept of "skewness" to measure the unbalance utilization of a server is introduced. By minimizing skewness, the overall utilization of servers in the face of multi-dimensional resource constraints can be improved.
- A load prediction algorithm is designed that can capture the future resource usages of applications accurately without looking inside the VMs. The algorithm can capture the rising trend of resource usage patterns and help reduce the placement churn significantly.

II. RELATED WORK

**Resource allocation by live VM migration**

A widely used technique for dynamic resource allocation in a virtualized environment is the VM live migration [3]. Sandpiper is also another mechanism which combines multi-dimensional load information into one single Volume metric. It arranges the list of PMs based on their volumes and the VMs in each PM in their volume-to-size ratio (VSR). This unfortunately diverts away critical information needed when making the migration decision. So it then considers the VMs and the PMs in the unsorted order.

The HARMONY system uses the virtualization technology across multiple resource layers. It also uses VM and data migration to mitigate hot spots not only on the servers, but also on the storage nodes as well as network devices. This system also introduces the Extended Vector Product (EVP) which indicates the imbalance use of resources. The load balancing algorithm used in harmony is another variant of the Toyoda method for multi-dimensional knapsack problem. But this system does not support green computing and load prediction. To minimize SLA violations, the dynamic placement of virtual servers is studied. Then they model it as a bin packing problem and the well-known first-fit approximation algorithm is used to calculate the VM to PM layout periodically. However, this algorithm cannot be used on-line; they are mostly designed for off-line use.

**Green Computing**

Many researches have been made to reduce energy consumption in data centers. Novel thermal design for lower cooling power, or adopting power-proportional and low-power hardware are hardware based approaches. Dynamic Voltage and Frequency Scaling (DVFS) is used to adjust CPU power consumption according to its load. For green computing, the DVFS is not used. PowerNap [7], a new hardware technologies that implements rapid transition (less than 1ms) between full operation and low power state such as Self-Refresh DRAM and Solid State Disk(SSD), so that it can “take a nap” in short idle intervals. Somniloquy is also another mechanism like power nap. For example, when a server goes to sleep, it notifies an embedded CPU residing on a special designed NIC to act as the main operating system. It gives the illusion that the server is always online. This paper belongs to the category of low-cost pure-software solutions. Similar to Somniloquy, SleepServer, another similar mechanism initiates virtual machines on a dedicated server as guest, instead of depending on a special NIC. Another mechanism that does not use a delegate is the LiteGreen. Instead it migrate the desktop OS away so that the desktop can sleep. In this category, it is necessary that the desktop is virtualized with shared storage. Jettison introduces “partial VM migration”, another variance of live VM migration, that only migrate away necessary working set while leaving unused data behind.

III. SYSTEM ARCHITECTURE

The architecture of the system is presented in figure shown below. In this architecture there may be 'n' numbers of servers and each server is hosting virtual machines VM1 and VM2. Each VM encapsulates one or more applications such as Facebook, Google etc. Several numbers of clients say 'n' number of clients may request services to the servers.

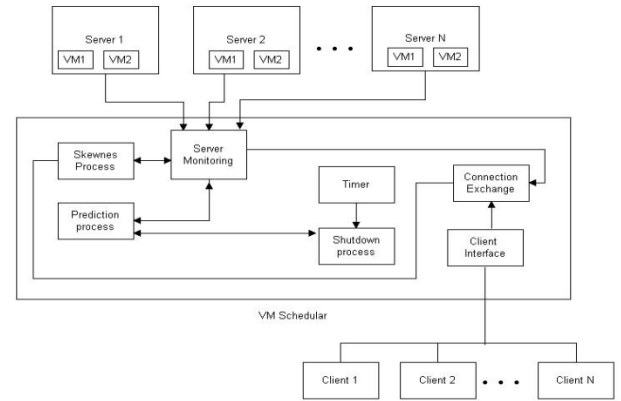


Figure: System Architecture

The Clients interact with the server through the client interface in the VM Scheduler. The Skewness Process calculates the HotSpot, ColdSpot and WarmSpot of the servers.

Hot spot:- If the utilization of any of its resources is above a hot threshold, then the server is in HotSpot. This shows that overloading of server resources occurs on the server and hence some VMs that run on the server can be migrated away to another server.

Cold spot:- If the utilizations of all its resources are below a cold threshold, then the server is in cold spot. This indicates that the server is mostly unuse and a potential candidate to shut down the server to save energy.

Warm spot:- We define a server as a warm spot if the utilizations of all its resources are below a hot threshold and below cold threshold. This indicates that the server is ready to run VMs.

The Server monitoring in the VM scheduler stores the information of the servers, for example the server monitoring has the information such as which server is experiencing more load i.e hotspot and which is experiencing less load i.e coldspot. And based on the information on the server monitoring and the number of clients on each server the prediction process predicts which server is needed to shut down. The server in Coldspot is then sent to the shut down process to shut it down and those clients in the mentioned server are sent to the Connection exchange. The shut down process also has a timer which periodically refreshed it and checks for incoming server for shutting down.

**Load Prediction Management**

The future resource need of Virtual machine (VMs) is needed to be predicted. Looking inside a VM for

application level statistics, e.g., by parsing logs of pending requests, is one particular solution.

To do this, the VM needs to be modified which may not always be possible. Instead, a prediction is made based on the past external behaviors of VMs.

In this paper, the exponentially weighted moving average (EWMA) prediction algorithm plays an important role in improving the stability and performance of our resource allocation decisions. Based on Physical machines (PMs) Usage we will select server using EWMA algorithm.

#### **Overload avoidance**

The existing capacity of PM should have sufficient resources to satisfy the resource needs of all VMs running on it. Else, the PM gets overloaded and may decrease the performance of its VMs.

The concept of "skewness" to measure the imbalance utilization of a server is introduced. By minimizing skewness, the overall utilization of servers can be improved in the face of multi-dimensional resource constraints. To perform skewness algorithm the following server status spot is required:

- **Hot spot**:- A server can be defined as a hot spot if the utilization of any of its resources is above a hot threshold. This indicates that the server is overloaded and hence some VMs running on it should be migrated away.
- **Cold spot**:- A server can be defined as a cold spot if the utilizations of all its resources are below a cold threshold. This indicates that the server is mostly idle and a potential candidate to turn off to save energy.
- **Warm spot**:- A server can be defined as a warm spot if the utilizations of all its resources are below a hot threshold and below cold threshold. This indicates that the server is ready to run VMs.

#### **Green computing**

The number of PMs used should be reduced as long as they can still satisfy the needs of all VMs. Idle PMs can be shut down to save energy consumption. When the server becomes cold spot, we should not give further more connections to that server and we will turn off the server. This will do by using Green computing.

#### IV. CONCLUSIONS AND FUTURE WORK

The Gateway system has intelligence to transfer the connection to cloud server in efficient way such that load will be distributed evenly in all cloud server.

In this project we have VM Scheduler and n number of server to perform the Dynamic resource allocation task. Server can get the connections from the client system. Based on the number of connection to the server, VM Scheduler will do the scheduling tasking (Migration, green computing). If number of connections exceeds the limit of server then VMscheduler involves in migration process. All incoming connection of server is migrated to another server via this same server. This is called Migration process.

Green computing is the task of closing/turnoff the devices which are not have current connections in the server.

Based on the number of connections we can perform the Migration and Green computing process.

#### REFERENCES

- [1] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield, "Xen and the art of virtualization," in Proc. of the ACM Symposium on Operating Systems Principles (SOSP'03), Oct. 2003.
- [2] K. Fraser, C. Clark, S. Hand, E. Jul, J. G. Hansen, C. Limpach, A. Warfield and I. Pratt, "Live migration of virtual machines," in Proc. of the Symposium on Networked Systems Design and Implementation (NSDI'05), May 2005.
- [3] M. Korupolu, A. Singh and D. Mohapatra, "Server-storage virtualization: integration and load balancing in data centers," in Proc. of the ACM/IEEE conference on Supercomputing, 2008.
- [4] S. Savage, Y. Agarwal and R. Gupta, "Sleepserver: a software-only approach for reducing the energy consumption of pcs within enterprise environments," in Proc. of the USENIX Annual Technical Conference, 2010.
- [5] V. N. Padmanabhan, T. Das, P. Padala, R. Ramjee, and K. G. Shin, "Litegreen: saving energy in networked desktops using virtualization," in Proc. of the USENIX Annual Technical Conference, 2010.
- [6] S. Hodges, R. Chandra, Y. Agarwal, R. Gupta and J. Scott, P. Bahl, "Somniloquy: augmenting network interfaces to reduce pc energy usage," in Proc. of the USENIX symposium on Networked systems design and implementation (NSDI'09), 2009.
- [7] E. d. Lara, N. Bila, K. Joshi, M. Hiltunen, M. Satyanarayanan and H. A. Lagar-Cavilla "Jettison: Efficient idle desktop consolidation with partial vm migration," in Proc. of the ACM European conference on Computer systems (EuroSys'12), 2012.
- [8] A. Vahdat, D. Gupta, G. M. Voelker and M. McNett, "Usher: Anvextensible framework for managing clusters of virtual machines," in Proc. of the Large Installation System Administration Conference (LISA'07), Nov. 2007.
- [9] B.-H. Lim, G. Hutchins and M. Nelson, "Fast transparent migration for virtual machines," in Proc. of the USENIX Annual Technical Conference, 2005.
- [10] K. Schwan and R. Nathuji "Virtualpower: Coordinated Power Management in Virtualized Enterprise Systems," Proc. ACM SIGOPS Symp. Operating Systems Principles (SOSP '07), 2007.