



## Energy Aware Resource Management of Cloud Data Centers

H. Rezai<sup>a</sup>, O. R. B. Speily<sup>\*b</sup>

<sup>a</sup> CoreBanking Research Group, Informatics Service Corporation, Tehran, Iran

<sup>b</sup> Department of Computer Engineering & Information Technology, Urmia University of Technology, Urmia, Iran

### PAPER INFO

#### Paper history:

Received 09 February 2017

Received in revised form 29 April 2017

Accepted 08 September 2017

#### Keywords:

Cloud Computing

Load Balancing

Power Saving

Virtualization

Live Migration

### ABSTRACT

Cloud Computing, the long-held dream of computing as a utility, has the potential to transform a large part of the IT industry, making software even more attractive as a service and shaping the way IT hardware is designed and purchased. Virtualization technology forms a key concept for new cloud computing architectures. The data centers are used to provide cloud services burdening a significant cost due to high energy consumption. Data centers are provisioned to accommodate peak demand rather than average demand and cloud applications consume much more electrical energy than they need. Thus, it necessitates that cloud computing solutions not only minimize operational costs, but also reduce the power consumption. In this paper, we investigate load balancing and power saving methods in virtualized cloud infrastructures. Imbalanced distribution of workloads across resources can lead to performance degradation and much electrical power consumption in such data centers. We present an architectural framework and principles for energy-efficient cloud computing environments. Resource provisioning and allocation algorithms, named Load-Power-aware, are proposed in this architecture. The algorithm employs a heuristic to dynamically improve the energy efficiency in data center, while guarantees the Quality of Service (QoS). The efficiency of the proposed approach is evaluated by using the most common cloud computing simulation toolkit, CloudSim. The performance modeling and simulation results are depicted the proposed approach significantly improves the energy efficiency in a given dynamic scenario, while a small amount of service level agreements (SLA) is missed.

doi: 10.5829/ije.2017.30.11b.14

## 1. INTRODUCTION

Today, there is a significant interest in developing data centers, in which applications are loosely coupled to the underlying infrastructure and can utilize the shared resources and It can be described as the supply of on-demand scalable resources as 'services' on 'pay-as-you-go' basis [1]. It is also desirable to facilitate the migration of an application from one set of resources to another without disruption in service. This is an important specification of the modern cloud computing infrastructures [2] that aims to efficiently share and manage extremely large data centers among applications. Two technologies that play an important role in cloud environments are virtualization [3] and live migration [4]. Virtualization provides an isolated environment to hosted virtual machines which have different operating systems and possess different

amount of resources on a shared physical platform [5]. This technology leads to virtualized physical hardware resources of data center and share them among multiple applications. The capabilities of virtualization are used to exploit the data center infrastructure for cloud computing and prepare flexible and reliable services. The processes and data of applications are hosted inside virtual containers (e.g., virtual machines and virtual storages) can be decoupled from the underlying physical resources. System virtualization is commonly used in infrastructure layer of cloud computing architecture. System virtualization enables different guest operating systems, virtual machines (VMs) to be concurrently executed on a bare machine, named physical host. The virtualizer software, called Virtual Machine Monitor (VMM), intercepts requests from VMs and conveys them to the corresponding physical resources. It also manages and coordinates virtual resources and the way underlying physical resources to be multiplexed among VMs. Xen [3], VMware ESX Server [6] are examples of

\*Corresponding Author's Email: speily@uut.ac.ir (O. R. B. Speily)

system virtualizers. Another technology used in cloud computing data center is live migration. This technology allows VMs and their contents to be migrated across distinct physical hosts without disruption. These technologies facilitate fault management, load balancing, power saving, and system maintenance in computing environments such as cloud computing (see [4] for more information). Load balancing in distributed systems highly improves system performance and increases resource utilization [7]. Load balancing is a process of reassigning the total load to the individual hosts of the collective system in order to affect resource utilization and response time improvement of submitted jobs [8]. Resource management aims to fairly distribute total loads of data center among resources and avoids hotspots. Overloaded hosts (hotspots) often lead to performance degradation and are vulnerable to failures [9]. To alleviate such hotspots, loads must be migrated to the underutilized hosts. In cloud computing environment, these loads are served inside VMs. When a host becomes overloaded, in fact, some VMs of that host must be migrated to appropriate hosts. There are two types of load balancing algorithms: Static load balancing; in this approach, priori knowledge about the global status of the distributed system, job resource requirement, and communication time are assumed [8]. A general disadvantage of all static schemes is that the final selection of a host to process allocated jobs is performed when the process is created and the allocation cannot be changed during process execution [10]. Some static algorithms are explained in [11, 12]. Dynamic load balancing; in this approach, load balancing decisions are based on the current state of the system. Tasks are allowed to move dynamically from an overloaded host to an under-loaded host to receive improved service. Although, these approaches are more complicated than static approaches, these approaches lead to a better performance. Because they make load balancing decisions based on the current system load [8]. Some dynamic algorithms are explained in [13, 14]. The growing consumption and decreasing access to natural resources result in increase in the price of energy [15]. Also, the rapid growth in demand for services that offered in data centers has led to consume large amount of electrical power which causes enormous cost burden to the cloud service providers. The biggest part of energy consumption is concerned with computation and transmission [16]. This cost includes hosts energy consumption and cooling and energy supply [17, 18]. According to Gartner report, 12 percent of overall data center cost relates to energy consumption [19]. As reported in [17], the greatest portion of data center cost relates to servers power consumption including CPU, memory, and storage systems. The aim of cloud data centers is to provide high performance while meeting SLAs and resource allocation is done to achieve this

purpose. But, without minimizing energy consumption, the most cost and energy will be wasted. By dynamic resource allocation methods, ensuring SLAs brings challenges to application performance management in virtualized environments. For this reason, cloud computing data centers must be designed to achieve not only the efficient processing and utilization of a computing infrastructure, but also to save and minimize energy usage [30]. Otherwise, cloud resources must be allocated not only to satisfy QoS requirements specified by users via SLAs, but also to reduce energy consumption. To achieve both performance and power efficiency is required effective consolidation policies that can minimize energy consumption without compromising the user-specified QoS requirements. In order to compare the efficiency of the proposed approach with others, we use several metrics to evaluate their performance. The first metric is the total energy consumption by the physical resources of a data center caused by the application workloads. The second metric called the average SLA is defined as the amount of SLA that could provide the level of QoS by provider based on the negotiated QoS requirements between the provider and consumers. The violated SLA, could be assumed that the provider pays a penalty to the clients in SLA violation (the third metric defined as penalty due to SLA violation). In this work, we investigate the cloud data center power consumption and develop an approach for power saving in cloud data center while satisfying SLA. There are several ways for decreasing power consumption of data centers including Dynamic Voltage and Frequency Scaling (DVFS) [20], virtualization of computer resources, and turning the idle servers off in order to further save on energy costs [21]. Our aim in this work is: 1) to describe the proposed architecture and principles for energy efficient cloud computing data center, 2) to investigate energy-aware resource placement and migration algorithms to improve the energy consumption of a data center, without violating the negotiated SLAs, and 3) to describe energy aware approach for managing of resources of data center and mapping VMs to suitable resources to use energy effectively. The rest of the paper is organized as follows: in Section 2, we give briefly discuss related work. The architecture and principles for management of clouds are presented in section 3. The proposed architecture is based on vSphere topology [22]. Placement and migration processes is explained in this section and the calculation model of power consumption and SLA calculation model is presented in the following parts of this section. In section 4, the proposed load balancing and power saving approach is described and the placement and migration algorithms applied in the proposed architecture is presented in more detail. In section 5, the proposed cloud data center and the proposed approach is evaluated by the cloud computing

simulation toolkit, CloudSim [23] and finally, Section 6 concludes the paper with summary and plans for future works.

## 2. RELATED WORKS

The studies [17, 18] indicate that the costs associated with the power consumption, cooling requirements, etc., of servers over their lifetime are significant and [18] shows that data centers are a significant source of CO<sub>2</sub> emissions. For this, power consumption of data centers has been studied in recent years. As said before, the use of virtual machine and migration technologies facilitates power saving in cloud computing environments, but improper use of these technologies, imposes additional overheads to the cloud data centers and consequently, leads to undesirable results in power saving and response times and affects the efficiency of cloud service. Hence, the use of optimum threshold evaluation methods for evaluation of hosts states and minimization of VM migrations reduces processing time and power consuming of data centers. In [24], an evaluation method has been proposed to dynamically evaluate the load status of the resources. It is a self-adaptive method based on load status of resources. In the study, first, some dynamic evaluation indicators has been presented to evaluate the resource load status more accurately. Second, the presented method divides the resources to three states, including Overload, Normal and Idle using the self-adaptive utilization threshold. Finally, an energy evaluation model has been presented to explain energy consumption using this method and the presented method has been compared to other methods. In [25], a methodology has been proposed to predict the availability of data center resources in future. The prediction is done by using genetically weight optimized artificial neural network and then the management of VMs are performed based on the availability of data center resources in future. In [26], an adaptive task scheduling strategy has been developed to achieve an optimal balance between energy efficiency ratio to performance in cloud tasks. In [27], coordinated power management has been studied in large-scale virtualized data centers. This is the first study that investigates the power management of virtual systems. The authors have presented a new method for power management, called soft resource scaling, by enumerating the benefits of virtualization regardless of the hardware scalability and virtual machines integration. The goals of the proposed approach in [27] are to online power management and to support the isolated and independent operations of guest virtual machines running on virtualized platforms. Authors in [28], applied a local search unit to minimize the power consumption and the cost of virtual machines migration.

The local search unit considers a utilization threshold level and if the server utilization is above the threshold level, a set of candidate virtual machines are selected for migration. Then it sorts all virtual machines based on processor utilizations and chooses the server with efficient power consumption as a destination. This method may result in too much immigration which causes the system thrashing.

## 3. CLOUD ARCHITECTURE MODEL

Our proposed architecture is based on vSphere topology [22]. This architecture comprises of several structural blocks including clusters of computational servers/hosts, networks and storage arrays, and local management unit that manages computational hosts and other resources. The infrastructure of the data center is composed of several heterogeneous hosts that virtual machines are created on them and run users applications. Every host characterized with given amounts of CPU, RAM and network bandwidth. Users submit requests starting the number of needed heterogeneous VMs for which the resource required MIPS, amount of RAM and network bandwidths are clarified. The local management unit (local manager) calculates amount of required resources to guarantee the user's SLA and minimize power consumption. The local manager powers off extra hosts for power management purposes in data center. Then, local manager chooses the most suitable host and directs user workload into it. If enough resources are available, then resources are instantly allocated to the request without waiting. Otherwise, on the occasion of insufficient hardware resources, VMs requests are put into a queue by the cloud queue unit. Coming requests wait in the queue for provisioning required resources. If a user requests more than one VM, local management unit may select one host for all VM requests or different hosts for each VMs based on host's CPU utilization. Local manager monitors the utilization level of hosts and specifies overloaded and under-loaded hosts based on a given threshold of utilization levels which are defined during hosts configuration. All of VMs running on the under-loaded hosts along with some VMs of overloaded hosts will volunteer to migrate lively to appropriate hosts. The live migration process of candidate VMs will be done based on the proposed Load-Power Aware algorithm, explained in the next section. As a result of the proposed method, migration of different VMs causes load balancing and power saving. Live migration of whole workloads of under-loaded hosts and powering them off leads to power saving in large scale cloud computing data centers. Also, live migration of several workloads submitted on overloaded hosts to the low utilized ones leads to load balancing and more power saving. This method is more

detailed in the following sections.

### 3. 1. VM Allocation and Migration Processes

Allocating resources in a cloud data center should be carried out carefully to minimize overall energy consumption while providing high performance and meeting SLA. We use the following relations to determine the most appropriate host during VM allocation and migration processes to measure the overloading degree of a host, we use the notion of an imbalance score, an exponential function, in Equation (1). The IBscore will be used in Equation (3) to select the best hosts for hosting VMs.

$$IBscore_i(u_i, T_{i,upper}) = e^{\alpha(u_i - T_{i,upper})} \quad (1)$$

where  $u_i$  is the utilization level of resource  $i$ , ranged from 0 to 1, of a given host over a resource such as CPU, RAM, storage, network interface, and etc.  $T_{i,upper}$  is the upper bound of utilization threshold of a given host over resource  $i$ . If the utilization level of resource  $i$  is smaller than  $T_{i,upper}$ , the value of  $(u_i - T_{i,upper})$  will be negative, and so its imbalance score factor may be smaller and comparison among hosts will not be correctly accomplished. The  $\alpha$  parameter is used to improve the correctness of IBscore. We need to calculate the portion of VM resources requirement on each candidate host. We define Load Fraction<sub>i,vm,host</sub> parameter which shows the fraction of needed resource  $i$  by VM  $v$  on host  $h$  as Equation (2).

$$LoadFraction_{i,vm,host} = \frac{\beta \times S_{i,r}}{S_{i,h}} \quad (2)$$

where,  $S_{i,r}$  is the amount of requested resource  $r$  by VM <sub>$i$</sub> ,  $S_{i,h}$  is the residual amount of resource  $i$  of host  $h$ , and  $\beta$  is a constant that initialized at the implementation time. As the fraction of resources requirements on each host impacts directly on load balancing and power saving in data center, we define the product of IBscore and Load Fraction as EVP<sub>i,host</sub> parameter as equation:

$$EVP_{i,host} = IBscore_{i,host} \times LoadFraction_{i,vm,host} \quad (3)$$

As the hosts are heterogeneous in a data center and their resources specification are different, the EVP<sub>i,host</sub> parameter accurately yields the best host selection for the submitted workloads. For example, consider two hosts A and B with 10000 and 3000 MIPS CPU capability, respectively, and a VM which needs 200 MIPS of a CPU to process a workload. Also, consider host A and B are utilized 60% and 55%, respectively, and  $\alpha=10$ ,  $\beta=10$ , and Tupper=80%. The EVP<sub>cpu</sub> for host A and B are approximately 0.027 and 0.064, respectively. It shows that the workload of VM should be hosted on host A.

For the best target host selection, the host with the lowest EVP parameter gains higher priority. Because

the suitable host for hosting the VMs has smaller IBscore, with lower utilization level, and smaller LoadFraction parameters. This means that if a VM created or migrated to that host, the probability that the host become overloaded is low and avoids thrashing occurrence.

To select an appropriate host, the EVP of each host is calculated based on all types of hardware resources and the dot product of IBscore and LoadFraction vectors. Each dimension of the vectors is related to one of host's resource types, for example CPU, RAM, network interface, etc. The dot product of the vectors determines attractiveness of each host to host specific VM based on different resources and requirements. Equations (4), (5), and (6) show that these parameters include all resource types.

$$IBscore_{host} = [IBscore_{s_1,host}, \dots, IBscore_{s_n,host}] \quad (4)$$

$$LoadFraction_{vm,host} = [LoadFraction_{s_1,vm,host}, \dots, LoadFraction_{s_n,vm,host}] \quad (5)$$

$$EVP_{host} = IBscore_{host} \cdot LoadFraction_{vm,host}^T \quad (6)$$

### 3. 2. Power Consumption Model

Power consumption of servers in data centers concerns utilization of CPU, memory, disk storage, and network interfaces. CPU is the main consumer of power in comparison with other resources in a server [28]. There are several techniques to save server power such as Dynamic Frequency Scaling (DFS), Dynamic Voltage Scaling (DVS), or a combination of Dynamic Voltage and Frequency Scaling (DVFS) that have been studied in recent studies [18, 28, 29] and shown that the application of these techniques on CPU causes a non-linear relationship between power and frequency at different CPU utilization levels. Moreover, the CPU utilization is typically proportional to the overall system load [29]. In [18, 22], a relationship between server CPU utilization and its power consumption has been presented. The idea is that power consumption increases by increasing the CPU utilization from idle state to fully utilized state. Moreover, these studies have shown that an idle server on average consumes approximately 70% of the power consumed by the server running at the full CPU speed [29]. The non-linear relationship between server power consumption and CPU utilization formulated as follows [18]:

$$P = P_{idle} + (P_{busy} - P_{idle}) \times (2u - u^r) \quad (7)$$

where,  $P_{idle}$  and  $P_{busy}$  are server power consumption in idle and busy states respectively. In this nonlinear model CPU utilization shown as  $u$  and  $r$  is a calibration parameter that minimize squared error of gathered data from the experimental system. Equation (7) shows that resource utilization is one of the influential factors of

power consumption in a data center. Although, data centers are provisioned to serve requests at peak, if the resources are underutilized, higher costs will be imposed to the cloud service providers[30]. A system with less than 50 percent resources utilization is virtually ineffective and has high overhead costs of maintaining and cooling systems. This necessitates a tradeoff between system performance and power costs. All servers in our model are assumed to follow this relationship. Additionally, the servers may be operated in an inactive mode by switching idle servers to the sleep mode to reduce more total power consumption and costs. This can be advantageous if the workload and utilization is low. Indeed, the minimum power  $P_{idle}$  is required to maintain the server in the active state which is typically substantial.

**3. 3. Service Level Agreement** One of the important requirements for a cloud computing environments is to provide reliable QoS. It is defined in terms of service level agreement (SLA) that describes such characteristics as minimal throughput, maximal response time or latency delivered by the deployed system. In this work, we calculate specific resource SLA violation for a VM as the following:

$$\overline{SLA}(t) = \frac{\sum_{\forall app \in VM} (requestedS_i(t) - allocatedS_i(t))}{\sum_{\forall app \in VM} requestedS_i(t)} \quad (8)$$

Average SLA violation for a host will be as the following:

$$\overline{SLA}_{host} = \frac{1}{n} \sum_{\forall vm \in host} \overline{SLA}_{vm} \quad (9)$$

where n is the number of created VMs on a host. If there exist m hosts in a data center, specific resource SLA violation at time t for that data center will be calculated as the following:

$$\overline{SLA}_{DataCenter} = \frac{1}{m} \sum_{\forall host} \overline{SLA}_{host} \quad (10)$$

and total average SLA violation at time t is obtained as the following:

$$ave\overline{SLA}_{DataCenter}(t) = \frac{1}{t} \int_{\tau=0}^t \overline{SLA}_{DataCenter}(\tau) \quad (11)$$

Finally, the average SLA data center can be obtained as:

$$aveSLA(t) = 1 - ave\overline{SLA}_{DataCenter}(t) \quad (12)$$

We use the last equation to demonstrate the portion of QoS that lost regarding to power saving and load balancing processes in a data center.

**3. 4. Penalty due to SLA Violation** There are significant advantages in the use of penalties due to SLA violation for both user and service providers. If a user receives all of its requested requirements from

cloud service provider and the service does not miss its deadline, the service is offered according to the SLA agreed with user. If the service misses its deadline, the user does not need to pay the whole price of received service out of required deadline. The penalty amount that cloud service provider pays the users, depends on the contract negotiated between them. In this situation, the cloud service provider pays a penalty due to SLA violation. The value of SLA violation is calculated according to the penalty function ( $\phi(t_v)$ ). This function relates to the delay time (td).  $t_v$  is the time that must be spent to complete the service after the deadline violation.  $t_v$  is calculated as follows(Equation (13)):

$$t_v = \overline{SLA} * t_r \quad (13)$$

where  $t_r$  is the response time of a service and  $\overline{SLA}$  is the SLA violation. Equation (14) shows response time of a service which includes  $t_v$  plus deadline  $t_d$ .

$$t_r = t_v + t_d \quad (14)$$

Response time is the time duration which an application is submitted to a data center until the service is completed. This relationship between penalty function and delay time may be linear or nonlinear. We assume that this relation is calculated by a linear function with a rate k. Then, the penalty due to SLA violation is given by Equation (15) as follows:

$$\phi(t_v) = price_{service} \times k \times t_v = price_{service} \times k \times \overline{SLA} \times t_r \quad (15)$$

where the amount of penalty that paid for violating the SLA within the  $t_r$ . So, the total cost for service providing in a cloud data center, using power-aware provisioning methods for power saving, will be calculated as the following:

$$C(t_r) = \phi(t_v) + C_p(t_d) \quad (16)$$

$$C(t_r) = [price_{service} \times k \times \overline{SLA} + price_{power} \times C_p(t_r)] \times t_r \quad (17)$$

where  $C(t_r)$  is the total cost imposed to cloud service provider within time  $t_r$  and  $price_{service}$  is the price per hour that user pays to cloud data center for receiving service. Also,  $price_{power}$  is the price per kWh paid for power consumption by cloud service providers.

#### 4. LOAD BALANCING & POWER SAVING APPROACH

When VMs are unfairly distributed across the different hosts and provided resources do not efficiently utilized by VMs, so the system may encounters hotspots and vulnerabilities.

Furthermore, it may impose more electrical power consumption and costs to the system. Data center power

consumption could be effectively reduced by supporting the migration of VMs among physical hosts and logically resizing the VMs to be consolidated within the minimum number of physical hosts. In addition, switching off the idle hosts can alleviate total consuming energy in a data center. For this purpose, we presented a dynamic power and load controlling method, named Load-Power-Aware, described as follows. In the dynamic methods, the current system state is considered for monitoring and controlling the system functionality [8].

---

#### Algorithm 1 VM Placement algorithm

---

VmAllocationPolicy Algorithm

input: VMList, output: allocation of Vm

minEVP ← Max

allocatedHost ← null

foreach (vm in VMList) do

  foreach (host in power-on-HostsList) do

    if (host has enough resource for vm) then

      EVP ← estimateEVP(host, vm)

      If (EVP < minEVP) then

        minEVP ← EVP

        allocatedHost ← host

      end-if

    end-if

  end-foreach

  If (allocatedHost = null) then

    foreach (host in power-off-Hosts) do

      if (host has enough resource for vm) then

        EVP ← estimateEVP(host, vm)

        If (EVP < minEVP) then

          minEVP ← EVP

          allocatedHost ← host

        end-if

      end-if

    end-loop

  end-if

  if (allocatedHost ≠ null) then

    allocate vm to allocatedHost

  end-if

end-loop

return allocation

---

The proposed method is aimed to avoid hotspots and distribute the submitted workloads fairly among resources in data center while saving power and guaranteeing SLAs. The idea is to allocate all VMs to hosts which any host utilization will be kept between lower and upper threshold levels. Thus, the overall data center load will be optimally balanced among hosts. If the CPU utilization of a host falls below the lower threshold, all VMs have to be migrated from the host and the host has to be powered off or switched to the

sleep mode in order to reduce the power consumption. If the utilization exceeds the upper threshold, some VMs have to be migrated from the host to reduce the utilization. The primary placement and the placement of migrated VMs is done based on our proposed placement algorithm. The pseudo-code for the placement algorithm is presented in Algorithm1. The algorithm assigns VMs to destination host based on the best fit scheme as it applies EVP parameter to find a host with the smallest product of IBscore and LoadFraction values, according to the model described in Section 3.1. The proposed method is carried out in two steps. At the first step, VMs needing to be migrated are chosen and at the second step, the selected VMs are placed on the hosts using the placement algorithm. The migrated VMs are chosen from under loaded and overloaded hosts.

---

#### Algorithm 2 VM migration algorithm

---

Input: hostList, VmList output: migrationList

getHostsStatus(hostList)

foreach h in hostList do

  vmList ← h.getVmList()

  vmList.sortIncreasingUtilization()

  hUtil ← h.getUtilization

    If (h is underLoaded) then

      migrationList.add (vmList)

    else if (h is overLoaded) then

      while (hUtil > h.getUpperUtilizationThreshold)

        foreach (vm in vmList) do

          migrationList.add(vm)

          hUtil ← hUtil – vm.getUtilization

        end-foreach

      end-while

    end-if

  end-foreach

return migrationList

---

The pseudo-code of VMs choosing algorithm for migration is presented in Algorithm2. The idea is to allocate all VMs to hosts which any host utilization will be kept between lower and upper threshold levels. Thus, the overall data center load will be optimally balanced among hosts. If the CPU utilization of a host falls below the lower threshold, all VMs have to be migrated from the host and the host has to be powered off or switched to the sleep mode in order to reduce the power consumption. If the utilization exceeds the upper threshold, some VMs have to be migrated from the host to reduce the utilization. Placement of migrated VMs is done based on our proposed placement algorithm which described earlier. The migration is done lively. Lively migration of a virtual machine makes the virtual machines operate without interruption. This migration type has low overhead. After completion of VM migration, running of VM on source host will be

terminated and allocated resources will be released upon residing on the target host. Then, migrated VM continues execution on the target host. Migration of VMs from overloaded hosts are performed periodically at interval time  $T$  while migration of VMs belong to underloaded hosts, to powering host off, is executed periodically at interval time  $nT$ ,  $n$  is a positive integer value which is set during implementation. As the main mission of cloud data center is service offering to customers then enough resources must be available every time to be allocated to submitted requests. Thus, the load balancing algorithm must be performed at short time duration. If powering off the hosts performs in a short interval times, then may be resources will not be available upon VMs submitting or migration phase. Furthermore, considering a short time interval may cause data center cost to increase as the overhead of powering on/off hosts occurs more times.

## 5. PERFORMANCE EVALUATION

Because of difficulties in running large-scale experiments on real-world infrastructure, simulation tools for evaluating will be a suitable way to evaluate the proposed method [25]. For evaluation of the proposed algorithm, the most common cloud computing simulation toolkit, CloudSim [25] is used in the reported work. The CloudSim toolkit supports modelling and creation of one or more virtual machines (VMs) on a simulated node of a Data Center, jobs, and their mapping to suitable VMs. We extend this framework in order to save the overall power and balance the load in cloud data center. The overhead of execution of the proposed algorithm is shown as SLA violation. Simulated cloud data center consists of 100 heterogeneous hosts which have virtualizable resources. Each host is modeled to have one core with 1000, 2000 or 3000 MIPS, 8 Gb of RAM, and 1 TB of storage. Every request consists of an application which runs inside a VM. Power consumption of each host is defined according to the model described in Section 3.2. According to this model, the power consumption of each host varies from 175 W (the idle state host) up to 250 W (fully utilized with 100% utilization). Also, the cost of power consumption is calculated according to the model described in Section 3.4. The price of energy consumption per kWh is assumed \$0.0677. The users submit requests for provisioning of 300 heterogeneous VMs. These requests fill the full capacity of the simulated data center. Each VM requires one CPU core with 250, 500, 750 or 1000 MIPS, 128 MB of RAM, and 1 GB of storage. Each VM runs an application with variable workload, which is modeled to generate the utilization of CPU according to a uniformly distributed random variable. We have simulated the applications as

EC2 Standard instance types. These instance types are small (default), medium, large, and extra-large. The small, medium, large and extra-large applications consist of maximum of 150,000 million, 750,000 million, 1,500,000 million, and 3,000,000 million instructions respectively. The unit price per hour for these applications equals to \$0.065, \$0.130, \$0.260, and \$0.520, respectively. The proposed algorithm evaluated in comparison with three power saving methods named Non-policy, only DVFS, and Single threshold (ST) [29]. Non-Policy method does not apply any optimization scheme to reduce the power consumption during run time and each host consumes maximum power all the time [29]. DVFS technology does not optimize the VM allocation at runtime, but can tune the host power consumption according to workload variations. We consider three power saving scenarios in cloud data center. In the first scenario, hosts in cloud data center are assumed to apply DVFS for power saving. In the second scenario, Single Threshold, hosts are able to use DVFS and upper threshold level to migrate VMs. In the third scenario, our proposed method, hosts apply DVFS for saving the power and powering host on/off. Also, it enhances the Load-Power-Aware algorithm to optimize VMs allocation and to migrate processes. Single Threshold method is based on the idea of setting the upper utilization threshold for hosts. The VMs are placed on hosts so that the total utilization remains below the upper threshold. If the host utilization exceeds the upper threshold, the candidate VM (VMs) will be lively migrated. At each time frame, the VMs are selected according to the VM migration algorithm. The selected VMs are reallocated to the new hosts according to the VM Placement algorithm. This causes overloaded hosts utilization reduction to the upper utilization threshold. In this experiment, upper utilization threshold of hosts is set to 80%. The ST method causes significant power consumption reduction in comparison to DVFS and Non-policy, but its profit is not better than DVFS due to SLA violation. In our proposed scheme, Load-Power-Aware algorithm is used to save further electrical energy and to fairly distribute the workloads across the hosts. For this purpose, this algorithm tries to keep the utilization of hosts between lower and upper threshold levels. This causes preventing hot spots and reducing energy consumption. Hosts target utilization is set to 60% and target utilization tolerance is set to of 20% (lower utilization threshold of 40% and upper utilization threshold of 80%). The results of comparison are depicted in Figures 1 to 4. The comparison of methods is evaluated for small, medium, large, and extra-large applications. Figure 1 depicts the energy consumption of the applications. As shown in Figure 1, the energy consumption of cloud data center by our proposed method is significantly reduced. VMs can only migrate

in the ST and Load-Power-Aware methods, as depicted in Figure 2, which causes the SLA violations. Total cost due to the SLA violation and power consumption in data center is shown in Figure 3. As shown in Figure 1, the total cost of ST method is further than DVFS method, this is because of SLA violations. Also, as depicted in Figure 4 the profit of ST method is lower than DVFS. But a significant profit has been achieved by our proposed method in comparison with other methods. We also evaluate our methods to determine the best target utilization and interval between lower and upper thresholds. The Load-Power-aware method is evaluated in three target utilization levels and host utilization interval which varies from 0.3 to 1. The results presented in Figures 5 show that decreasing the energy consumption accompanies with increasing of SLA violations. Figure 6 depicts the profit achieved by the Load-Power-Aware method. The profits from different target utilization and interval between upper and lower thresholds show that the best combination of the target utilization and interval between the thresholds is achieved when the target utilization and the interval between the thresholds are set to 60% and 40%, respectively.

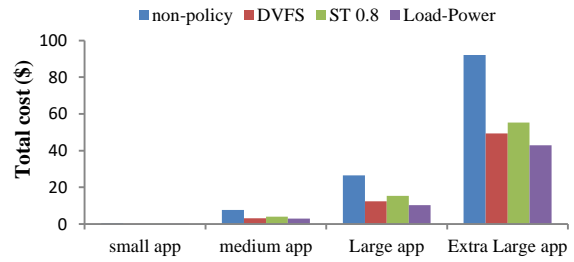


Figure 3. Total cost SLA violation and power consumption

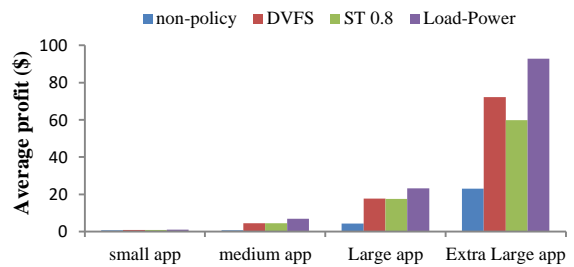


Figure 4. The average Profit achieved by different methods

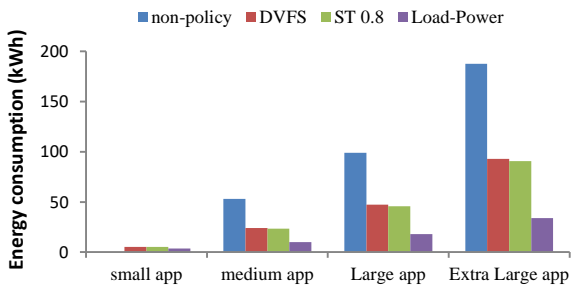


Figure 1. Energy consumption by different methods

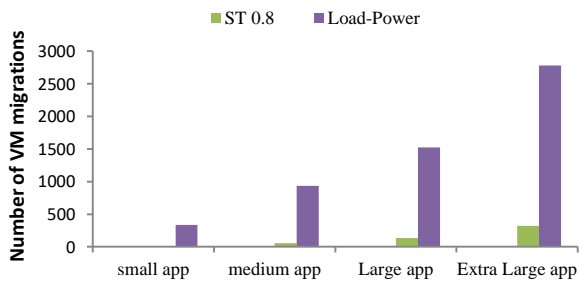


Figure 2. The number of VM migrations

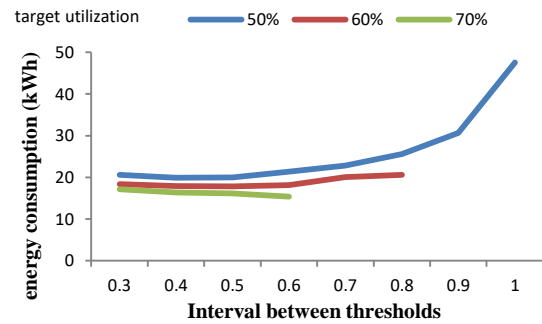


Figure 5. The energy consumption

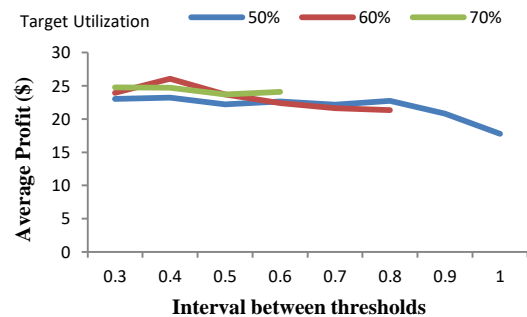


Figure 6. The average profit



## 6. CONCLUSIONS

In this work, the problem of dynamic resource allocation and power management in cloud computing data centers is investigated. We have presented Load - Power-aware resource allocation algorithm that utilizes dynamic consolidation of VMs to reduce the power consumption and enhance the data center profit. The experiment results have shown that this approach leads to a substantial reduction of energy consumption in cloud data centers and significant achieved profit from the cost of power consumption perspective, in comparison to other resource allocation techniques. The proposed method is able to maintain host utilization between specified intervals and power off idle hosts to achieve more electrical energy saving. Moreover, it prevents hotspots in data center. As a part of future works, we plan to enhance the proposed algorithm to optimize other resource types and to consider the algorithm for network optimization. In order to reduce data transfer overhead and network devices load, it is crucial to consider network communication between VMs in real location decisions. So development of the proposed algorithm for the network optimization is a subject of our future research works. We also plan to study load balancing and power saving in federated clouds [33] and evaluate our approach for these distributed data centers.

## 7. REFERENCES

- Jeyanthi, N., Shabeeb, H., Durai, M.S. and Thandeeswaran, R., "Rescue: Reputation based service for cloud user environment", *International Journal of Engineering-Transactions B: Applications*, Vol. 27, No. 8, (2014), 1179-1186.
- Buyya, R., Yeo, C.S., Venugopal, S., Broberg, J. and Brandic, I., "Cloud computing and emerging it platforms: Vision, hype, and reality for delivering computing as the 5th utility", *Future Generation Computer Systems*, Vol. 25, No. 6, (2009), 599-616.
- Barham, P., Dragovic, B., Fraser, K., Hand, S., Harris, T., Ho, A., Neugebauer, R., Pratt, I. and Warfield, A., "Xen and the art of virtualization", in ACM SIGOPS operating systems review, ACM. Vol. 37, (2003), 164-177.
- Clark, C., Fraser, K., Hand, S., Hansen, J.G., Jul, E., Limpach, C., Pratt, I. and Warfield, A., "Live migration of virtual machines", in Proceedings of the 2nd Conference on Symposium on Networked Systems Design & Implementation-Volume 2, USENIX Association., (2005), 273-286.
- Smith, J. and Nair, R., "Virtual machines: Versatile platforms for systems and processes, Elsevier, (2005).
- Devine, S.W., Bugnion, E. and Rosenblum, M., "Virtualization system including a virtual machine monitor for a computer with a segmented architecture". (2002), Google Patents.
- Khan, Z., Singh, R., Alam, J. and Kumar, R., "Performance analysis of dynamic load balancing techniques for parallel and distributed systems", *International Journal of Computer and Network Security*, Vol. 2, No. 2, (2010), 123-127.
- Alakeel, A.M., "A guide to dynamic load balancing in distributed computer systems", *International Journal of Computer Science and Information Security*, Vol. 10, No. 6, (2010), 153-160.
- Singh, A., Korupolu, M. and Mohapatra, D., "Server-storage virtualization: Integration and load balancing in data centers", in Proceedings of the 2008 ACM/IEEE conference on Supercomputing, IEEE Press., (2008), 53-60.
- Sharma, E., Singh, S. and Sharma, M., "M.: Performance analysis of load balancing algorithms", in In: 38th World Academy of Science, Engineering and Technology, Citeseer., (2008).
- Tang, X. and Chanson, S.T., "Optimizing static job scheduling in a network of heterogeneous computers", in Parallel Processing, 2000. Proceedings. 2000 International Conference on, IEEE., (2000), 373-382.
- Motwani, R. and Raghavan, P., "Randomized algorithms", *ACM Computing Surveys (CSUR)*, Vol. 28, No. 1, (1996), 33-37.
- Karimi, A., Zarafshan, F., Jantan, A., Ramli, A.R. and Saripan, M., "A new fuzzy approach for dynamic load balancing algorithm", *arXiv preprint arXiv:0910.0317*, (2009).
- Zeng, Z. and Veeravalli, B., "Rate-based and queue-based dynamic load balancing algorithms in distributed systems", in Parallel and Distributed Systems, 2004. ICPADS 2004. Proceedings. Tenth International Conference on, IEEE., (2004), 349-356.
- Rouholamini, M. and Mohammadian, M., "Grid-price-dependent energy management of a building supplied by a multisource system integrated with hydrogen", *International Journal of Engineering-Transactions A: Basics*, Vol. 29, No. 1, (2016), 40-49.
- Taghizadeh, S. and Mohammadi, S., "Lebrp-a lightweight and energy balancing routing protocol for energy-constrained wireless AD HOC networks", *International Journal of Engineering-Transactions A: Basics*, Vol. 27, No. 1, (2013), 33-40.
- Greenberg, A., Hamilton, J., Maltz, D.A. and Patel, P., "The cost of a cloud: Research problems in data center networks", *ACM SIGCOMM computer Communication Review*, Vol. 39, No. 1, (2008), 68-73.
- Fan, X., Weber, W.-D. and Barroso, L.A., "Power provisioning for a warehouse-sized computer", in ACM SIGARCH Computer Architecture News, ACM. Vol. 35, (2007), 13-23.
- Gartner, I., Gartner says energy-related costs account for approximately 12 percent of overall data center expenditures. (2010), Tech. Rep.
- Semeraro, G., Magklis, G., Balasubramonian, R., Albonesi, D.H., Dwarkadas, S. and Scott, M.L., "Energy-efficient processor design using multiple clock domains with dynamic voltage and frequency scaling", in High-Performance Computer Architecture, 2002. Proceedings. Eighth International Symposium on, IEEE., (2002), 29-40.
- Urgaonkar, R., Kozat, U.C., Igarashi, K. and Neely, M.J., "Dynamic resource allocation and power management in virtualized data centers", in Network Operations and Management Symposium (NOMS), 2010 IEEE, (2010), 479-486.
- Introduction to vmware infrastructure: Esx server 3.5, esx server 3i version 3.5, virtualcenter 2.5. (2007), Revision.
- Buyya, R., Ranjan, R. and Calheiros, R.N., "Modeling and simulation of scalable cloud computing environments and the cloudsim toolkit: Challenges and opportunities", in High Performance Computing & Simulation, 2009. HPCS'09. International Conference on, IEEE., (2009), 1-11.

24. Zuo, L., Shu, L., Dong, S., Zhu, C. and Zhou, Z., "Dynamically weighted load evaluation method based on self-adaptive threshold in cloud computing", *Mobile Networks and Applications*, Vol. 22, No. 1, (2017), 4-18.
25. Radhakrishnan, A. and Kavitha, V., "Energy conservation in cloud data centers by minimizing virtual machines migration through artificial neural network", *Computing*, Vol. 98, No. 11, (2016), 1185-1202.
26. Shen, Y., Bao, Z., Qin, X. and Shen, J., "Adaptive task scheduling strategy in cloud: When energy consumption meets performance guarantee", *World Wide Web*, Vol. 20, No. 2, (2017), 155-173.
27. Nathuji, R. and Schwan, K., "Virtualpower: Coordinated power management in virtualized enterprise systems", in *ACM SIGOPS Operating Systems Review*, ACM. Vol. 41, (2007), 265-278.
28. Verma, A., Ahuja, P. and Neogi, A., "Pmapper: Power and migration cost aware application placement in virtualized systems", in *Proceedings of the 9th ACM/IFIP/USENIX International Conference on Middleware*, Springer-Verlag New York, Inc. (2008), 243-264.
29. Beloglazov, A., Abawajy, J. and Buyya, R., "Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing", *Future Generation Computer Systems*, Vol. 28, No. 5, (2012), 755-768.

## Energy Aware Resource Management of Cloud Data Centers

H. Rezaei<sup>a</sup>, O. R. B. Speily<sup>b</sup>

<sup>a</sup> CoreBanking Research Group, Informatics Service Corporation, Tehran, Iran

<sup>b</sup> Department of Computer Engineering & Information Technology, Urmia University of Technology, Urmia, Iran

### PAPER INFO

چکیده

#### Paper history:

Received 09 February 2017

Received in revised form 29 April 2017

Accepted 08 September 2017

#### Keywords:

Cloud Computing

Load Balancing

Power Saving

Virtualization

Live Migration

رایانش ابری، به عنوان یک رویای بلندمدت و ابزار رایانشی، این قابلیت را دارد که بخش بزرگی از صنعت فناوری اطلاعات را تغییر داده و صنعت نرم افزار را به عنوان سرویس ارائه شده جذابتر سازد و شکل طراحی و فروش سخت افزار لازم برای فناوری اطلاعات را تغییر دهد. فناوری مجازی سازی مفهوم کلیدی معماری های جدید رایانش ابری را تشکیل می دهد. مراکز داده ای که برای ارائه سرویس های ابری مورد استفاده قرار می گیرند، هزینه سنگینی را به خاطر مصرف انرژی بالا تحمیل می کنند و نرم افزارهای کاربردی ابری، بیش از حد مورد نیازشان انرژی الکتریکی مصرف می کنند. لذا، روش های رایانش ابری نه تنها باید هزینه های عملیاتی را کاهش دهند، بلکه باید مصرف توان را نیز کاهش دهند. در این مقاله، روش های تعدیل بار و صرفه جویی مصرف توان در زیرساخت های مجازی سازی شده ابری را بررسی کرده ایم. توزیع نامتعادل بارهای کاری بین منابع، منجر به کاهش بازدهی و افزایش مصرف توان الکتریکی در چنین مراکز داده می شود. ما، یک چارچوب معماری و اصولی برای محیط رایانش ابری با کارایی انرژی بالا ارائه کرده ایم. الگوریتم فراهم سازی و تخصیص منابع، که Load-Power-aware نام گذاری شده، در این معماری پیشنهاد شده است. این الگوریتم، در ضمن حفظ کیفیت سرویس، یک روش ابتکاری پویا را که مصرف انرژی را بهبود می بخشد، به کار می برد. کارایی روش ارائه شده با ابزار شبیه سازی رایانش ابری، CloudSim ارزیابی شده است. نتایج مدل سازی و شبیه سازی کارایی، نشان می دهند که روش ارائه شده کارایی مصرف توان را به صورت قابل توجهی بهبود می بخشد در حالی که مقدار ناچیزی از توافق سطح ارائه سرویس از بین می رود.

doi: 10.5829/ije.2017.30.11b.14