

# Effects of Correlation-based VM Allocation Criteria to Cloud Data Centers

Jing V. Wang\*, Chi-Tsun Cheng, and Chi K. Tse  
 Department of Electronic and Information Engineering  
 The Hong Kong Polytechnic University  
 Hunghom, Kowloon, Hong Kong  
 \*Email: jing.j.wang@connect.polyu.hk

**Abstract**—Virtualization technology has been widely adopted in Cloud data centers for adaptive resource provisioning. With virtualization, multiple virtual machines (VMs) can be co-located on a single physical host to yield maximum efficiency. However, VMs which show high CPU utilization correlations to other co-located peers are more likely to trigger overloading incidents. This work provides an analysis on effects of correlation-based VM allocation criteria to Cloud data centers. The correlations among VMs' CPU utilizations are considered as parameters for decision making in VM allocation processes. Three different expressions of correlation-based criteria are introduced and evaluated in this work. According to our simulation results obtained from CloudSim with real-world workload traces, Cloud data centers with correlation-based allocation criteria can perform better in terms of reducing energy consumption and avoid committing Service Level Agreements violations than those with power-based criteria.

**Keywords**—Resource Provisioning, Cloud Computing, VM Allocation, Correlation, CPU utilization

## I. INTRODUCTION

Cloud technology, by pooling resource to on-demand computing in a cost-effective manner, is gaining prominence rapidly. The soaring demand for Cloud applications has produced a surge in resource utilization of Cloud data centers. Resource provisioning, which minimizes the number of active physical hosts by allocating virtual machines (VMs) carefully, is an efficient way to reduce energy expenditure of Cloud data centers. On the other hand, it is essential for Cloud service providers to provide the committed processing power to their subscribers, or a penalty cost will be applied. Virtualization technology allows resource of a physical host to be shared by multiple VMs. However, hosts with highly correlated VMs are more likely to trigger overloading incidents. Therefore, how to prevent the co-location of highly correlated VMs on the same host becomes an important issue that needs to be addressed.

Resource allocation problems in Cloud computing can be regarded as combinatorial problems. Several techniques have been proposed to analyze performance interference effects between co-located VMs. In [1], Zhu and Tung proposed a consolidation algorithm based on an interference model to search an optimal consolidation configuration. Nathuji *et al.* in [2] proposed a QoS-aware control framework to manage performance interference effects introduced by the

consolidation of multiple VMs onto multicore servers. Their work is based on a MIMO model to determine whether additional resources should be allocated to compensate performance degradation due to interference between co-located workloads. However, both of these works were focusing on interference effects in VM consolidation processes and did not emphasis on VM migration techniques. In [3], a correlation-aware dynamic power management solution targeting the execution of scale-out applications was presented by Kim *et al.*. They considered the correlations between co-located VMs individually instead of calculating the multiple correlation. In an earlier work of the authors in this paper [4], host's temperature was used as a migration criterion. In [5], the provisioning process was formulated as a stable matching problem to make hosts operate at desirable utilization levels.

In this work, several VM allocation criteria based on VMs' CPU utilization correlations are presented and analyzed. The criteria were utilized in an allocation mechanism which optimizes resource allocation to mitigate overloading caused by correlated VMs. More specifically, correlation information among VMs are taken into account in the migration process to lower the risk of further overloading on source hosts while without imposing negative impacts on destination hosts. The criteria were implemented and evaluated on CloudSim [6] with real-world workload data. Comparing with allocation mechanisms with power-based criteria, mechanisms with correlation-based criteria show significant improvements in terms of energy consumption and fulfilling Service Level Agreements (SLAs).

The rest of the paper is arranged as follows. Preliminaries are given in Section II. Section III introduces and elaborates the correlation-based VM allocation criteria. Performance of the correlation-based mechanism is studied and discussed in Section IV. Finally, conclusions are given in Section V.

## II. PRELIMINARIES

We adopt the multiple correlation coefficient in [7] to estimate the correlation between VMs. Consider a host with  $n$  co-located VMs and suppose these VMs are represented by vector  $\mathbf{V} = [V_1, V_2, \dots, V_n]$ . The correlation strength of the  $i^{\text{th}}$  VM toward the other  $n-1$  VMs is measured based on their last  $q$  CPU utilization observations. Let  $\mathbf{y}_i$  be denoted as the vector containing the last  $q$  observations of the  $i^{\text{th}}$

VM. Similarly, let  $\mathbf{X}$  be denoted as an augmented matrix comprises the  $q$  observations of the remaining  $n - 1$  VMs on the host. Expressions of  $\mathbf{y}_i$  and  $\mathbf{X}$  are shown as follow

$$\mathbf{y}_i = \begin{bmatrix} y_{1,i} \\ \vdots \\ y_{q,i} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,m} & \cdots & x_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{p,1} & \cdots & x_{p,m} & \cdots & x_{p,n-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{q,1} & \cdots & x_{q,m} & \cdots & x_{q,n-1} \end{bmatrix}. \quad (1)$$

Here,  $\mathbf{y}_i$  contains the last  $q$  utilization history of the  $i^{\text{th}}$  VM, while  $\mathbf{X}$  contains that of all other co-located VMs. Here,  $x_{p,m}$  is the  $p^{\text{th}}$  CPU utilization observation of  $V_m$ . Then we can compute the multiple correlation coefficient  $R_{V_i, \mathbf{V} \setminus V_i}^2$  for each  $V_i$ , which is denoted as

$$R_{V_i, \mathbf{V} \setminus V_i}^2 = \frac{\sum_{k=1}^q (y_{k,i} - m_{\mathbf{y}_i})^2 (y_{k,i} - m_{\hat{\mathbf{y}}_i})^2}{\sum_{k=1}^q (y_{k,i} - m_{\mathbf{y}_i})^2 \sum_{k=1}^q (y_{k,i} - m_{\hat{\mathbf{y}}_i})^2}, \quad (2)$$

where  $m_{\mathbf{y}_i}$  and  $m_{\hat{\mathbf{y}}_i}$  are the sample means of  $\mathbf{y}_i$  and  $\hat{\mathbf{y}}_i$  respectively, and  $\hat{\mathbf{y}}_i$  is a vector of predicted values which can be obtained as follow

$$\hat{\mathbf{y}}_i = \mathbf{X}\mathbf{b} \quad \mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_i. \quad (3)$$

The correlation coefficient between the  $i^{\text{th}}$  VM and all the other co-located VMs is then estimated accordingly.

### III. CORRELATION-BASED VM ALLOCATION CRITERIA

In general, a Cloud resource provisioning process contains three major steps, which are (1) identifying over or under-utilized hosts, (2) selecting VM(s) on the identified hosts for migration, and (3) reallocating those VM(s) according to some given criteria.

In the first step, we adopt Local Regression Robust (LRR) algorithm introduced in [8] to identify overloaded hosts because of its superior performance comparing with other host overloading detection methods. LRR method is an adaptive utilization threshold detection method. It estimates the CPU utilization of a host based on its last  $j$  CPU utilization values, and thus determines whether a host is considered as overloaded. In this work,  $j$  is set to 10.

In the second part of the provisioning process, Minimum Migration Time (MMT) policy in [8] is adopted for VM selection. Under MMT, a VM associated with the shortest migration time on a critical host will be selected to be migrated first. In the results presented in [8], it is shown that MMT outperforms other selection policies in the VM selection step.

In this section, we introduce and elaborate three different correlation-based VM allocation criteria which can be utilized in the last step of the provisioning process, where suitable hosts will be identified to accommodate the migrated VM(s). The VM reallocation process is commonly formulated as a Bin Packing Problem (BPP). Among the

solvers for BPP, the Best-Fit-Decreasing (BFD) heuristic is employed in this work due to its low complexity.

#### A. Correlation of Migrated VM(s)

In this approach, a VM will be allocated to a host such that the correlation between the migrated VM and the existing VMs on the host is minimized. Such correlation is calculated using (2).

#### B. Average Correlation Level of Destination Host(s)

In the second approach, we allocate a migrated VM to a host with the minimal average correlation level. A host's average correlation function (ACL) is defined as follow

$$\text{ACL} = \frac{\sum_{i=1}^n R_{V_i, \mathbf{V} \setminus V_i}^2}{n}, \quad (4)$$

where  $n$  is the total number of VMs on the candidate host together with the migrated in VM. Comparing with the previous approach, the current approach considers the impact of the migration to the co-located VMs and allows a host with a relatively large number of VMs being selected, provided that the correlations between the migrating-in VM and the co-located VMs are all at low values.

#### C. Variation of Correlation Level of Destination Host(s)

The higher the correlations among VMs running on a host, the higher the probability for the host to be overloaded [9]. Based on such phenomenon, in the last approach, we try to consolidate VMs such that each active host could achieve a low correlation level among all its co-located VMs to reduce the risk of overloading.

Intuitively, VMs with strong correlations should be placed onto different hosts to reduce such a risk. A VM to be migrated will not choose hosts with VMs that have strong correlations with it. It should also avoid causing significant performance impacts on the destination host. Therefore, we compute the total correlation variation of each candidate host to efficiently quantify the impact of the migrating-in VM(s) on their existing VMs, which is defined as

$$\text{VCL} = \sum_{i=1}^{n-1} \left( R_{V_i, \mathbf{V} \setminus V_i}^2 - R_{V_i, \mathbf{V}' \setminus V_i}^2 \right), \quad (5)$$

where  $\mathbf{V}'$  is the vector represented VMs on the host before receiving the migrating-in VM. Under this criterion, we select hosts with minimal VCL values for VM reallocation. All the above approaches can be applied in BFD algorithm for solving the re-allocation problem. Due to space limitation, only the last approach is presented in Figure 1.

Details on the operation of the mechanism are elaborated as follows. At first, we initialize a list of available hosts from the host overloading detection process and a list of to-be-migrated VMs (VmsToMigrate) obtained from the VM selection process. Then the selected VMs are sorted in a descending order of their current CPU utilizations. For each

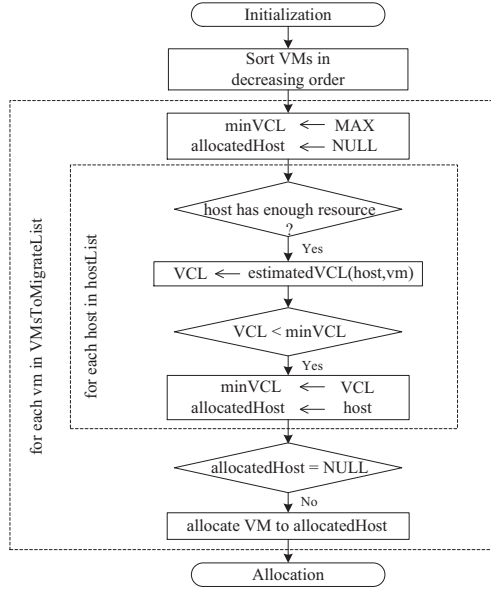


Figure 1: Flowchart of Correlation-based BFD Algorithm

VM in the pipeline, the host with the minimum VCL value will be selected as its destination. After each reallocation, the migrated VM would be removed from the VmsToMigrate list. If no host is available, an inactive host will be turned on to accommodate the VM. On the other hand, under-utilized hosts will be turned off to conserve energy. The algorithm is repeated until all the VMs on the VmsToMigrate list are being allocated.

It can be observed that an allocation process with the correlation-based criteria tends to consolidate VMs onto some hosts that can minimize the impact to their other co-located VMs. By doing so, an overloading incident is less likely to be triggered and the number of migrations can be reduced ultimately. Furthermore, mechanisms with these criteria can arrange VMs with low correlations to operate under the same host to yield better utilization.

#### IV. EXPERIMENTAL EVALUATION

##### A. Simulation Setup

We implemented and evaluated all the aforementioned criteria on CloudSim-3.0.3 [6]. In this work, the scenario under study is based on an Infrastructure as a Service (IaaS) model. An ordinary Cloud data center with 800 physical hosts and 1052 VMs was simulated. In the simulation, there were two types of dual-core hosts with different resource capacities, namely HP ProLiant G4 servers and HP ProLiant G5 servers. The corresponding energy models were obtained from SpecPower08 [10]. At any given instance, multiple independent users may submit their requests on provisioning  $M$  VMs. These VMs, characterized by their requirements, are then allocated to the physical hosts. To simulate real-world scenarios, four different types of single-core VMs with

different levels of MIPS and RAM were simulated in the simulations. Each VM was configured with 100 Mbit/s of bandwidth and 2.5 Gigabytes of storage. In the experiments, the sampling interval of overloading measurements is set to five minutes. The CPU utilization history  $q$  equals to 30. During the provisioning process, SLAs, a measurement of QoS, will be established between the Cloud service provider and its users. If there are any SLA violations, the service provider will have to pay a penalty, which will increase its operating cost.

In this paper, the level of SLA violation is measured using the two metrics in [8] : (1) SLA violation Time per Active Host (SLATAH) which indicates the percentage of time when physical hosts have reached 100% CPU utilization, and (2) Performance Degradation due to Migrations (PDM) which shows the overall performance degradation due to the capacity requirement of the migrated VM and the VM migration process itself. SLATAH and PDM are independent to each other and are with equal importance. A parameter called SLA Violation (SLAV), which integrated both metrics, is defined as

$$\text{SLAV} = \text{SLATAH} \times \text{PDM}. \quad (6)$$

In general, energy consumption and SLAV are conflicting metrics. The goal of a VM allocation mechanism is to achieve a reasonable trade-off between these two metrics. In this work, the Energy and SLA Violations (ESV) in [8] is adopted, that combines energy consumption and SLAV metrics together to evaluate the overall performance of a Cloud data center. Here, ESV is expressed as

$$\text{ESV} = E \times \text{SLAV}, \quad (7)$$

where  $E$  is the total energy consumption of a data center.

##### B. Performance Analysis

In our experiments, we chose the power-based LRR mechanism in [8] as a benchmark due to its outstanding performance over other existing methods. Here, the power-based LRR mechanism is referred to the method which adopted host's power consumption as a migration criteria. While in other mechanisms under test, correlation-based criteria mentioned in Section III were adopted separately.

Six different metrics were used to evaluate the efficiency of the correlation-based criteria, namely energy consumption, migration number, total overloaded hosts, SLATAH, SLAV, and ESV. Table I shows the results obtained from the simulations.

As observed in Table I, mechanisms utilizing the correlation-based criteria invoked less migrations compared to the benchmark. Note that the number of overloaded hosts obtained using correlation-based allocation mechanisms are much lower than that of the power-based mechanism. Table I also shows the SLATAH values of Cloud data centers with different VM allocation criteria. Correlation-based allocation

Table I: Simulation Results

VM Allocation Criteria	Energy (kWh)	Migration Number	Overloaded Hosts	SLATAH(%)	SLAV (x0.00001)	ESV (x0.001)
Power-based	163.48	27859	3138	5.86	4.64	7.59
Correlation of Migrated VM(s)	125.2	9792	2953	4.28	1.1	1.34
Average Correlation Level (ACL)	125.46	10147	3008	4.6	1.26	1.58
Variation of Correlation Level (VCL)	124.59	9546	2796	3.68	0.89	1.11

mechanisms led to significantly less SLATAH than their counterparts, which indicates that correlation-based allocation mechanisms can further reduce the risk of overloading and thus have less impact on the quality of service. Systems with lower ESV values mean they can achieve a better all-round performance. From the simulations, it can be observed that correlation-based mechanisms can outperform the power-based LRR mechanism by about 80% in terms of ESV. The simulation results show the advantages of considering correlation information among VMs during the VM reallocation process.

From the results, we find that allocation mechanisms adopting the VCL criterion can yield lowest ESV values among all other variations under test. Mechanisms utilizing the VCL criterion can allocate VM to hosts that introduce the least impacts to their co-located VMs. Furthermore, by minimizing the correlations among co-located VMs, the probability of having overloading incidents is reduced and thus lead to a low number of VM migrations. In contrast, for power-based method, hosts with different hardware configurations may encounter the same utilization level, but associated with very different energy consumption, and vice versa. Using energy consumption or utilization as the sole allocation criterion may lead to a poor resource utilization.

Note that the criterion based on ACL did not perform well as the other two correlation-based criteria. Criterion based on ACL works well if all the correlations between the migrating-in VM and the co-located VMs are at low values and are evenly distributed. However, it cannot identify cases when there exist a few extreme correlation values as they are averaged out by the large number of VMs.

## V. CONCLUSIONS

In this work, we presented an analysis on the effects of different correlation-based virtual machine (VM) allocation criteria to Cloud data centers. Correlation-based allocation mechanisms allocate VMs to hosts based on CPU utilization correlations among VMs. VMs with low correlations are more preferred to be co-located on the same physical host to lower the risk of overloading, and thus avoid potential Service Level Agreement (SLA) violations. Performances of correlation-based allocation mechanisms were evaluated using CloudSim with real-world workload data. Simulation results show that the criterion considers correlation variation of candidate hosts during a VM reallocation process

performs better than other VM reallocation criteria in terms of reducing energy consumption and SLA violations.

## ACKNOWLEDGMENT

This work is supported by the Department of Electronic and Information Engineering, the Hong Kong Polytechnic University (Projects G-YBKX and RTMR).

## REFERENCES

- [1] Q. Zhu and T. Tung, "A performance interference model for managing consolidated workloads in QoS-aware clouds," in *Cloud Computing (CLOUD), 2012 IEEE 5th International Conference on*. IEEE, 2012, pp. 170–179.
- [2] R. Nathuji, A. Kansal, and A. Ghaffarkhah, "Q-clouds: managing performance interference effects for QoS-aware clouds," in *Proceedings of the 5th European conference on Computer systems*. ACM, 2010, pp. 237–250.
- [3] J. Kim, M. Ruggiero, D. Atienza, and M. Lederberger, "Correlation-aware virtual machine allocation for energy-efficient datacenters," in *Proceedings of the Conference on Design, Automation and Test in Europe*. EDA Consortium, 2013, pp. 1345–1350.
- [4] J. V. Wang, C.-T. Cheng, and C. K. Tse, "A power and thermal-aware virtual machine allocation mechanism for cloud data centers," in *Communication Workshop (ICCW), 2015 IEEE International Conference on*, June 2015, pp. 2850–2855.
- [5] J. V. Wang, K.-Y. Fok, C.-T. Cheng, and C. K. Tse, "A stable matching-based virtual machine allocation mechanism for cloud data centers," in *2016 IEEE World Congress on Services (SERVICES)*, June 2016, pp. 103–106.
- [6] R. N. Calheiros, R. Ranjan, A. Beloglazov, C. A. F. De Rose, and R. Buyya, "Cloudsim: A toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms," *Softw. Pract. Exper.*, vol. 41, no. 1, pp. 23–50, Jan. 2011.
- [7] H. Abdi, "Multiple correlation coefficient," *The University of Texas at Dallas*, 2007.
- [8] A. Beloglazov and R. Buyya, "Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in cloud data centers," *Concurr. Comput. : Pract. Exper.*, vol. 24, no. 13, pp. 1397–1420, Sep. 2012.
- [9] A. Verma, G. Dasgupta, T. K. Nayak, P. De, and R. Kothari, "Server workload analysis for power minimization using consolidation," in *Proceedings of the 2009 conference on USENIX Annual technical conference*. USENIX Association, 2009, pp. 28–28.
- [10] "Specpower08," (Accessed: 2015-10-29). [Online]. Available: <http://www.spec.org>