



# Heuristics based server consolidation with residual resource defragmentation in cloud data centers



K. Sunil Rao, P. Santhi Thilagam\*

Department of CSE, NITK Surathkal, India

## HIGHLIGHTS

- We model residual resource fragmentation in cloud data centers.
- We examine feasibility of residual resource defragmentation with server consolidation.
- Proposed approach performs defragmentation with low energy cost and SLA violations.
- Controlled defragmentation with consolidation reduces VM migrations.
- A good mix of cloud applications ensures better defragmentation.

## ARTICLE INFO

### Article history:

Received 28 April 2014

Received in revised form

15 September 2014

Accepted 19 September 2014

Available online 24 October 2014

### Keywords:

Server consolidation

Cloud data center

IaaS

Virtualization

Bin packing

Resource fragmentation

## ABSTRACT

Server Consolidation is one of the foremost concerns associated with the effective management of a Cloud Data Center as it has the potential to accomplish significant reduction in the overall cost and energy consumption. Most of the existing works on Server Consolidation have focused only on reducing the number of active physical servers (PMs) using Virtual Machine (VM) Live Migration. But, along with reducing the number of active PMs, if a consolidation approach reduces residual resource fragmentation, the residual resources can be efficiently used for new VM allocations, or VM reallocations, and some future migrations can also be reduced. None of the existing works have explicitly focused on reducing residual resource fragmentation along with reducing the number of active PMs to the best of our knowledge. We propose RFAware Server Consolidation, a heuristics based server consolidation approach which performs residual resource defragmentation along with reducing the number of active PMs in cloud data centers.

© 2014 Elsevier B.V. All rights reserved.

## 1. Introduction

With rapid rise in the popularity of Cloud Computing, cloud data centers have started hosting a wide range of varied user applications in their PMs [1]. These applications can either be placed in dedicated PMs or colocated with other applications. In order to account for application performance, infrequent and inevitable workload peaks and security requirements, hardware isolation can be provided to the applications with minimum sharing of resources [2]. But such hardware isolation leads to sub-optimal resource utilization. Deploying only one application in a PM would bloat up the number of active PMs to be managed by a cloud provider, and operating a cloud data center with this model would be practically impossible. Moreover, resource pooling is one of the

fundamental characteristics of cloud computing [3] and any cloud provider needs to incorporate resource pooling in order to run a data center in a profitable manner.

Therefore, to exploit the advantages of multi-tenancy without affecting the application performance and security requirements, server virtualization technology [4] is used. With virtualization of servers, more than one virtual server can co-exist in a single PM which increases the PM resource utilization and reduces overall cost and energy consumption. Even though virtualization technology aims at maximizing resource utilization, most of the existing cloud data centers have utilization varying from 20% to 50%. This problem is popularly known as VM Sprawl [5]. VM Sprawl is one of the biggest challenges being faced by many companies owning virtualized data centers. VM Sprawl is a consequence of quick and uncontrolled creation of VMs. This results in PMs being provisioned unnecessarily without proper judgement and over-provisioning of resources to VMs where VMs consume resources more than what they actually require. Moreover, dynamic variation in the resource requirements of the applications deployed in a cloud environment

\* Corresponding author.

E-mail addresses: [sunilrao91@gmail.com](mailto:sunilrao91@gmail.com) (K. Sunil Rao), [santhi@nitk.ac.in](mailto:santhi@nitk.ac.in) (P. Santhi Thilagam).

results in rapid variation in the resource utilization of PMs. Hence there is a need to arrange VMs in the PMs intelligently at regular intervals.

### 1.1. Server consolidation

Server Consolidation, in a broader perspective, refers to the process of cutting down the Total Cost of Operation (TCO) of a data center either by relocating PMs in order to reduce the number of data center locations, or by reducing the total number of PMs required using virtualization letting more than one VM run on a PM. Our work focuses on reducing the total number of active PMs using VM Live Migration which is also known as Physical Consolidation [6]. Server Consolidation can be either be done offline or online. The references to Server Consolidation in the remaining part of this paper actually refers to Offline Physical Consolidation unless explicitly stated otherwise.

The problem of minimizing the number of active PMs can be formalized as a Vector Packing problem, which is a variant of Bin Packing problem, well known in the field of Operations Research. The Vector Packing problem is an NP Complete problem [7] and hence finding an exact solution for any given input can be done only in exponential time. Different approximation approaches have been proposed in the state of the art for server consolidation using different techniques. Most of the heuristics based approaches such as Harmonic and Cardinality Constrained Harmonic approach [8], First Fit Decreasing and Best Fit Decreasing [9], Modified Best Fit Decreasing [10], Improvised First Fit Decreasing [6], Sercon [11] are applications of the most widely known bin packing algorithms like Next Fit, First Fit and Best Fit algorithms [8]. Approaches such as pMapper [2], Dynamic Round Robin Approach [12], Genetic Algorithm Based Approach [13], Adaptive Threshold based approach [14], 2-Phase Optimization Method [15], Control Theoretic solution [16] provide various models and techniques for reducing power and energy consumption of the data center. Along with reducing power and energy consumption, approaches such as LP formulation and heuristics based approach [17] try to control VM migration taking VM capacity into consideration.

Most of the approaches aim at migrating VMs from under-utilized PMs to other PMs so that the under-utilized PMs can be set to a power saving state. In order to decide on the best destination PM for any migration, most of the existing approaches rely on evaluating the scores of the PMs based on various factors such as resource utilization, power and energy consumption. The problem of server consolidation is a multi-objective optimization problem [18] consisting of objectives such as improving resource utilization, reducing energy consumption, reducing SLA violations, reducing number of Live Migrations and reducing Residual Resource Fragmentation. The approach of using a scalar value as a score to decide on the destination PM for a migration cannot optimize multiple objectives simultaneously. For example, considering the objective of improving resource utilization only, it is not possible to improve utilization of CPU, Memory and Network Bandwidth simultaneously using a single scalar score. If a scalar score is used to decide on the destination PM, it leads to fragmentation of residual resources. The problem of resource fragmentation renders the residual resources useless or less useful, thereby adding to the cost incurred to the data center provider.

None of the existing works have explicitly focused on reducing residual resource fragmentation to the best of our knowledge. We intend to exploit the fact that since server consolidation being considered is offline consolidation which runs at regular intervals, there exists a scope for improving VM allocations to reduce residual resource fragmentation. Moreover, since the problem of server consolidation is a multi-objective optimization problem, we propose multi-phase approach for server consolidation with each

phase focusing on an individual objective. The idea here is to ensure that even though it is not possible to optimize all the parameters completely, it is possible to ensure that each parameter can be optimized as per requirement and can be made to fall in a tolerable range by using multiple phases. Our work focuses on exploring heuristics which can reduce residual resource fragmentation to make residual resources more useful and reduce energy consumption and cost incurred to the data center.

The contributions of this paper are as follows.

1. The paper presents the problem of residual resource fragmentation in the context of server consolidation and its consequences.
2. It proposes a model for quantifying and evaluating residual resource fragmentation.
3. The paper proposes a heuristics based multi-phase approach for server consolidation which effectively reduces residual resource fragmentation along with reducing the number of active PMs.

### 1.2. Organization of the paper

The remaining part of the paper is organized as follows. Section 2 describes the problem of residual resource fragmentation and its consequences. Section 3 outlines the problem statement and objectives. Section 4 discusses the solution methodology. Section 5 describes the experimental setup and the workloads considered in the implementation. Section 6 discusses the results and analysis, and Section 7 provides the conclusion and the future scope of our work.

## 2. Residual resource fragmentation

### 2.1. Residual resource

Residual Resource refers to the free resource available in the active PMs of a data center. Since VMs have dynamically varying resource requirements, almost all active PMs have non zero residual resources.

### 2.2. Residual resource fragmentation

Residual Resource Fragmentation refers to the state of the data center where sufficient amount of residual resources are available for any new VM allocation or VM reallocation, but are fragmented and distributed across multiple active PMs, rendering them unusable. The process of reducing residual resource fragmentation is called Residual Resource Defragmentation.

Consider the scenario in Fig. 1(a) consisting of 3 active PMs and 9 VMs. The total residual CPU resource is  $15 + 20 + 35 = 70\%$ . This residual CPU resource is spread across three active PMs and the maximum CPU that can be allocated for any VM without performing any migration is 35%. Similarly the maximum Memory resource that can be allocated for any VM is 25%.

But, if the VMs are placed as in Fig. 1(b), the maximum CPU resource that can be allocated to any VM is 50% and the maximum Memory resource that can be allocated to any VM is 40%. This is because the residual resource fragmentation is less in the scenario of Fig. 1(b) compared to that in scenario of Fig. 1(a).

### 2.3. Advantages of residual resource defragmentation

The advantages of Residual Resource Defragmentation are as follows.

1. It improves the usability of the residual resources in the data center. That is, if the residual resources are concentrated in less

Download English Version:

<https://daneshyari.com/en/article/425625>

Download Persian Version:

<https://daneshyari.com/article/425625>

[Daneshyari.com](https://daneshyari.com)