

# Chapter 1

## Introduction

### 1.1 Problem Statement

Cloud computing is a computing model offering a network of servers to their clients in a on-demand fashion. From NIST's definition [28], *"cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction."* To illustrate how it works, considering a case: a Cloud provider builds a data center which contains thousands of servers connected with network. These servers are virtualized which means they are partitioned into a smaller unit of resources called *Virtual Machines (VMs)*. A web-based application provider can access and deploy their applications (e.g Endnote, Google Drive and etc.) in these VMs from anywhere in the world. Once the applications start serving, application users can use them without installing on their local computers.

One of the major characteristic of Cloud computing is to separates the role of traditional service provides into Cloud user (software provider) and Cloud (infrastructure) provider. As Wei [47] states, "one provides the computing of services, and the other provides the services of computing". Therefore, stakeholders of Cloud computing become: Cloud providers, Cloud users, and End (application) users [19] (see Figure 1.1). This separation beneficial for both Cloud user and End user: It releases the burden of purchasing and maintaining hardwares for Cloud users. Consequently, lower the expense of End users.

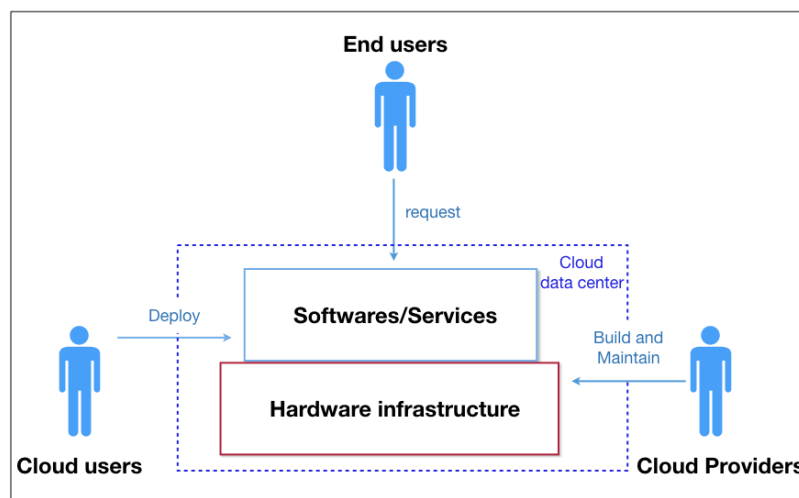


Figure 1.1: Stakeholders of Cloud computing

Each stakeholder has their responsibility, goal, and objectives.

- *Cloud providers* build data centers, provide maintenance and resource management on the hardware infrastructure. Their goal is to increase the profit by boosting the income and reducing the expense. Their income comes from Cloud users' rental of servers or *Physical Machines (PMs)* in terms of resource quality (e.g 3.5GHz dual-core CPU), resource quantity (e.g 3 PMs), and time (e.g 1 year). Therefore, Cloud providers objective is to maximize utilization of computing resources. A high utilization brings two benefits, firstly, it increases income by accommodating more applications in limited resources. Secondly, it cuts the expense of energy consumption by packing applications in a minimum number of servers so that idle servers can be turned off.
- *Cloud users* develop and deploy applications on Cloud. Each application generates time-vary CPU utilization. Their goal is also increase the profit mainly through two objectives, attracting more End users and reduce the expense of resources. The first objective can be achieved by improving the quality of service as well as lower the fee for End users. Either way depends not only on the software entities but also the quality of service (QoS) offered by Cloud provider. The second objective can be achieved by a good estimation of the reserved resources, so that they do not rent insufficient resources which cause low QoS or rent too much resources which cause wastage.
- *End Users* are the final customers in this chain. They consume services directly from Cloud users and indirectly from Cloud provider. Their goal is to obtain a satisfactory service. Their goal is achieved by signing a Service Level Agreement (SLA) with Cloud users which constrains the behavior of the services.

In this proposal, we focus on Cloud providers' perspective which is to reduce the energy consumption in Cloud data centers by providing a series of methods to improve the overall resource utilization.

Besides the stakeholders of Cloud computing, Cloud computing has three service models [28] (see Figure 1.2): Infrastructure as a Service (IaaS), Platform as a Service (PaaS) and Software as a Service (SaaS). Each service model describes distinct responsibilities of stakeholders and different resource management strategies.

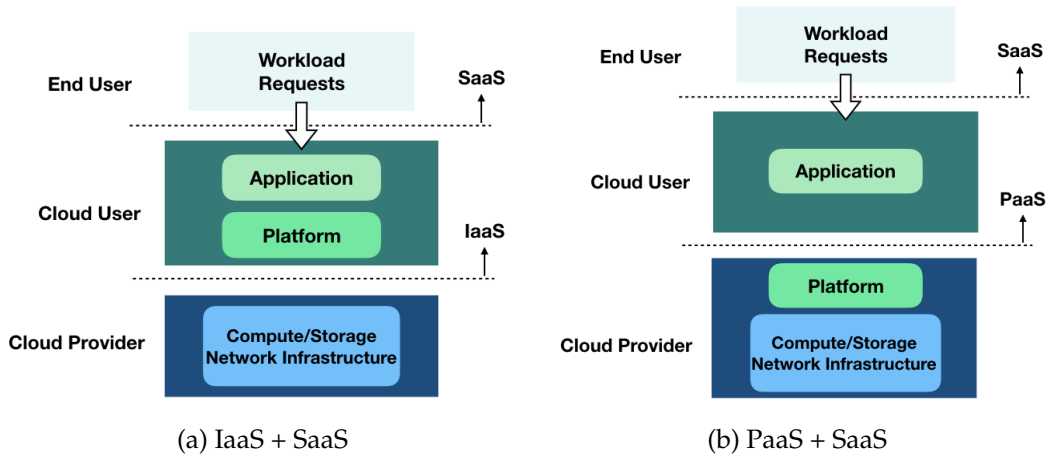


Figure 1.2: Three Service models [19]

- *IaaS*, Cloud provider offers the fundamental computing resources, often in the form of various sizes of VMs. Apart from the virtualized hardware and operating systems,

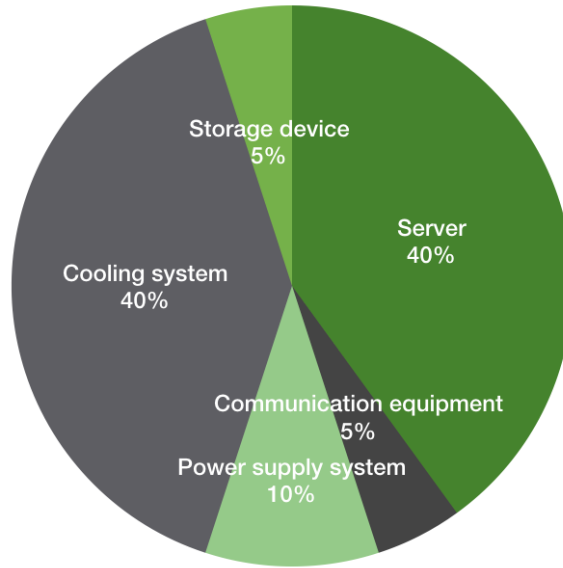


Figure 1.3: Energy consumption distribution of data centers [35]

Cloud users treat the remote VMs as local and deploy their applications. In terms of resource management, Cloud users have the responsibility to estimate the quantity of resources, while Cloud providers have no knowledge and control inside VMs. Resource management is based on VM.

- *PaaS*, Cloud providers establish the software development platform to enable the in-progress software to be developed in the platform. The main difference between PaaS and IaaS is that, in PaaS, Cloud providers not only offer the hardware but also a platform which allows software development. In terms of resource management, Cloud providers have the full control of resource allocation. Therefore, Cloud users can focus on software development.
- *SaaS* only describes the relationship between Cloud user and End User, therefore, it does not involve with resource management. Cloud users deploy their applications in Cloud which can be accessed by End Users.

Cloud computing has completely reformed the software industry [7] by providing three major benefits to Cloud users. First, Cloud users do not need upfront investment in hardware (e.g. PMs and networking devices) and pay for hardware's maintenance. Second, Cloud users do not need to worry about the limited resources which can obstruct the performance of their services when unexpected high demand occurs. The elastic nature of cloud can dynamically allocate and release resources for a service. In addition, Cloud users can pay as much as the resource under a *pay-as-you-go* policy. Third, Cloud users can publish and update their applications at any location as long as there is an Internet connection. These advantages allow anyone or organization to deploy their software on Cloud in a reasonable price.

Energy consumption [23] is the major concern of Cloud providers. It is derived from several parts as illustrated in Figure 1.3. Regardless of the energy consumption of the cooling system, the majority is from PMs. According to Hameed et al [16], PMs are far from energy-efficient. The main reason for the wastage is that the energy consumption of PMs remains high even when the utilization is low (see Figure 1.4). Therefore, a concept of *energy pro-*

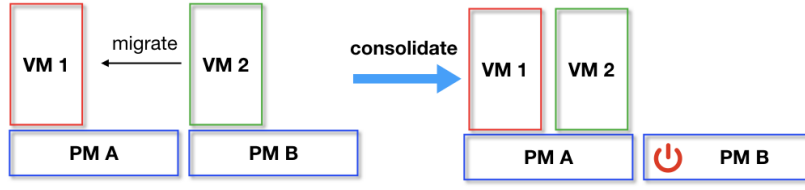


Figure 1.5: A Server Consolidation example: Initially, each PM runs an application wrapped with a VM in a low resource utilization state. After the consolidation, both VMs are running on PM A, so that PM B is turned off to save energy [3].

*portional computing* [3] raised to address the disproportionate between utilization and energy consumption.

The solution of the high energy consumption problem comes two important aspects: Virtualization technology which enables a fine granularity resource management and various strategies which are tight coupling with service models.

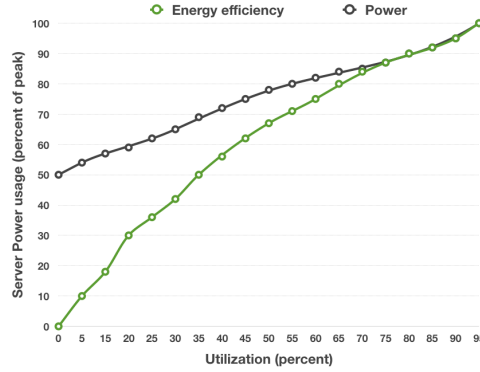


Figure 1.4: Disproportionate between utilization and energy consumption [3]

Virtualization [44] is the fundamental technology that enables Cloud computing. It partitions a physical machine's resources (e.g. CPU, memory and disk) into several isolated units called virtual machines (VMs) where each VM allows an operating system running on it. This technology rooted back in the 1960s' and was originally invented to enable isolated software testing.

In a Cloud data center, each server should not run arbitrary numbers of applications, because it cannot guarantee the performance of each application. Therefore, initially, each application is assigned to a server, which causes great wastage. VMs can provide good isolation which means applications running in co-located VMs within the same PM do not interfere each other [39]. Soon, people realized that it can be a way to improve the utilization of hardware resources: With each application deployed in a VM, a PM can run multiple applications.

Server consolidation [53] resolves the low utilization problem by gathering applications into a fewer number of PMs (see Figure 1.4), so that the resource utilization of PMs are maintained at a high level and the idle PMs can be turned off to save energy. Consolidation dramatically improves hardware utilization and lowers PM and cooling energy consumption.

Server consolidation is the core functionality involving in all Cloud resource manage-

ment operations. These operations can be roughly separated into three [42, 29] (see Figure 1.6) which are applied in different scenarios: Application initialization, Global consolidation, and Dynamic resource management.

1. *Application initialization* is applied when new applications or new VMs arrive and the problem is to allocate them into a minimum number of PMs.

In this problem, a set of applications or VMs are waiting in a queue. The resource capacity of the PM and usage by applications are characterized by a vector of resource utilizations including CPU, memory and etc. Then, the allocation system must select a minimum number of PMs to accommodate them so that after the allocation, the resource utilizations remain high. The problem is to consider the different combinations of applications so that the overall resource utilization is high. This problem is naturally modeled as a static bin-packing problem [10] which is a NP-hard problem meaning it is unlikely to find an optimal solution of a large problem.

2. *Global consolidation* is conducted periodically to adjust the current allocation of applications so that the overall utilization is improved.

In this problem, time is discrete and it can be split into basic time frames, for example: ten seconds. A periodical operation is conducted in every  $N$  time frames. A cloud data center has a highly dynamic environment with continuous arriving and releasing of applications. Releasing applications cause hollow in PMs; new arrivals cannot change the structure of current allocation. Therefore, after the initial allocation, the overall energy efficiency is likely to drop along with time elapsing. In global consolidation, an optimization system takes the current allocation - including a list of applications/VMs and a list of PMs - as the input, adjust its allocation so that the global resource utilization is improved.

In comparison with initialization, instead of new arrivals, the global consolidation considers the previous allocation. Another major difference is that global consolidation needs to minimize the differences of allocation before and after the optimization. This is because the adjustment of allocation relies on a technique called live migration [?], and it is a very expensive operation because it occupies the resources in both the host and the target. Therefore, global optimization must be considered as a time-dependent activity which makes the optimization even difficult.

3. *Dynamic resource management* Dynamic resource management is applied in three scenarios. **First**, it is applied when a PM is overloading. Overloading is defined as, one of the resources' consumption exceeds the PM's capacity. At this moment, the applications in that PM are suffering from a Quality of Service (QoS) degradation which means their performance cannot meet the agreement because of the limited resources. Overloading is often caused by the increasing of workload. In order to prevent the QoS from dropping, an application is migrated to another PM. This is called hot-spot mitigation [29]. **Second**, it is applied when a PM is under-loading. Under-loading is when a PM is in a low utilization state normally defined by a threshold. At this moment, all the applications in the under-loading PM are migrated to other active PMs, so the PM becomes empty and can be turned off. This is called dynamic consolidation. **Third**, it is applied when a PM having very high level of utilization while others having low. An adjustment is to migrate one or more application from high utilized PMs to low ones. This is called load balancing.

No matter which scenario it is, a dynamic resource management always involves three steps .

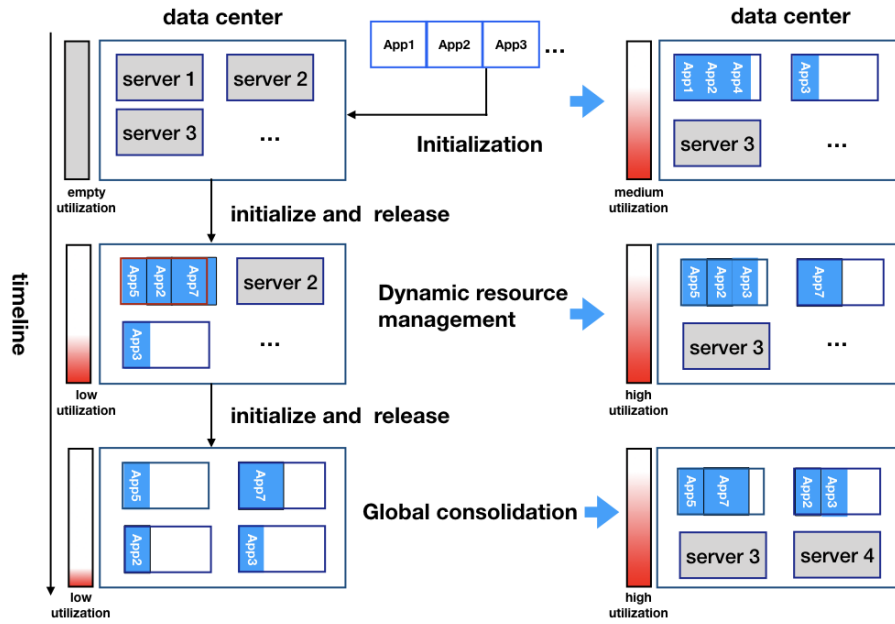


Figure 1.6: Cloud data center resource management involves three steps: initialization, dynamic resource management, and global consolidation. Grey means closed PMs.

- *When to migrate?* refers to determine the time point that a PM is overloaded or underloaded. It is often decide by a threshold of utilization.
- *Which application to migrate?* refers to determine which application need to be migrated so that it optimize the global energy consumption.
- *Where to migrate?* refers to determine which host that an application is migrated to. This step is called dynamic placement which is directly related to the consolidation, therefore, it is decisive in improving energy-efficiency.

Among three operations, dynamic placement is a dynamic and on-line problem. The term “dynamic” means the request comes at an arbitrary time point. An on-line problem is a problem which has on-line input and requires on-line output [?]. It is applied when a control system does not have the complete knowledge of future events.

There are two difficulties in this operation, firstly, dynamic placement requires a fast decision while the search space is very large (e.g hundreds of thousands of PMs). Secondly, migrate one application at a time is hard to reach a global optimized state.

Finally, a consolidation plan includes four major items:

1. A list of existing PMs after consolidation
2. A list of new virtual machines created after consolidation
3. A list of old PMs to be turned off after consolidation
4. The exact placement of applications and services

By the nature of Cloud resource management, server consolidation techniques can also be categories into static and dynamic methods [49, 45]. Static method is a time consuming

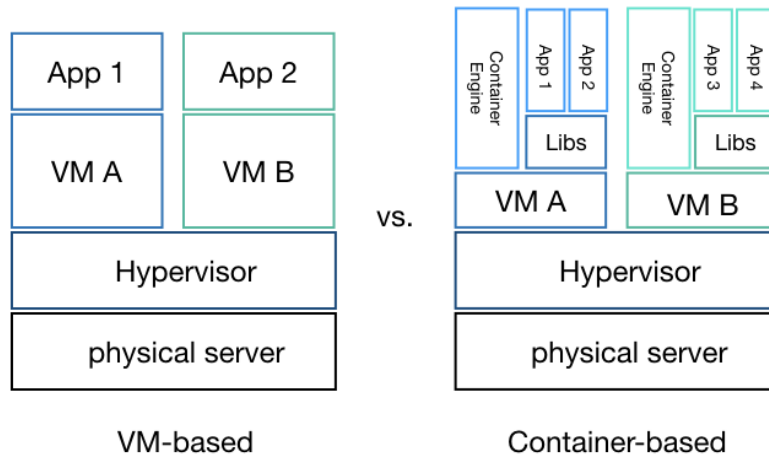


Figure 1.7: A comparison between VM-based and Container-based virtualization

process which is often conducted off-line in a periodical fashion; initialization and global consolidation belong to this category. It provides a global optimization to the data center. Dynamic method adjusts PMs in real time. It often allocates one application at a time. Therefore, it can be executed quickly and often provides a local optimization to the data center.

Resource management strategies are very different in IaaS and PaaS. In IaaS, Cloud users are responsible for the estimation of resources which is represented as fixed types of VMs and Cloud providers based on the size of VMs to consolidate them into fewer PMs. This strategy has two inherent flaws which lower the resource utilization.

- First one is the separated responsibilities of resource selection for applications and resource allocation. Because of separated responsibilities, Cloud users must estimate the quantity of resources. They tend to reserve more resources for ensuring the QoS at the peak hours [9]. This causes the low utilization.
- Second is the fixed types of VMs, where each type of VM represents a certain amount of resources (e.g. CPU, RAM, and Storage). Because of the fixed size of VM and the one-on-one mapping of applications and VMs (see Figure 1.7 left hand-side), specific applications consume unbalanced resources which leads to vast amount of resource wastage [43]. For example, computation intensive tasks consume much more CPU than RAM; a fixed type of VM provides much more RAM than it needs. Because the tasks use too much CPU, they prevent other tasks from co-allocating. This causes wastage.

Traditionally, *overbooking* strategy is used to improve the resource utilization in IaaS. Overbooking [43] or oversubscribe is strategy to allocate more VMs in a PM even their aggregated resources exceed the PM's capacity. This strategy aims at solving the problem of over-provisioned VMs. The advantage of this approach is it consolidates more VMs in PMs and, indeed, improved the resource utilization. The disadvantages are also obvious: An overbooked PM is easy to suffer from overloading which is its capacity cannot satisfied the required requests.

In recent years, virtualization technology has evolved to allow finer granularity resource management. A recent development of Container technique [38] has driven the attention

of both industrial and academia. Container is an operating system level of virtualization which means multiple containers can be installed in a same operating system (see Figure 1.7 right-hand side). Each container provides an isolated environment for an application. In short, a VM is partitioned into smaller manageable units. This new concept starts a new service model called Container as a Service (CaaS) [34]. CaaS brings advantages for both Cloud customers and providers. From Cloud users' perspective, CaaS has advantages of both IaaS (Infrastructure as a Service) and PaaS (Platform as a Service) but without their disadvantages. On one hand - similar to PaaS - it does not require Cloud users to estimate the quantity of resources so that they can focus on application development. On the other hand - similar to IaaS - it allows Cloud users to customize their software environment without being constrained by platforms.

For Cloud providers, CaaS resolves two IaaS's inherent weaknesses which cause low utilization of resources.

In contrast, CaaS solves above two problems at a time. It allows Cloud providers to manage both the resource selection for applications and resource allocation; It also enables VM-resizing and many-to-one mapping between applications to VMs (Figure 1.7 right hand-side). Hence, Cloud providers have a complete control of resources which may result in a better utilization of resources. In addition, IaaS Cloud runs many redundant operating systems and hypervisors. CaaS eliminates these redundancies by a single operating system with multiple containers.

Currently, vast amount of server consolidation methods are mostly VM-based and they are mainly modeled as bin-packing problems [25], where applications represent items and PMs represent bins. These methods can not be directly applied on container-based problem because container-based consolidation has two levels of allocation: Containers are allocated to VM as the first level and VM are allocated to PM as the second level. These two levels of allocation interact with each other.

Even for single level of bin-packing problems, the complexity is NP-hard. In VM-based static problem, deterministic methods such as Integer Linear Programming [41] and Mixed Integer Programming [46] are often considered. However, it is well-known that they are very time-consuming for a large scale problem. More research proposed heuristic methods to approximate the optimal solution such as First Fit Decreasing (FFD) [32], Best Fit Decreasing (BFD) [4]. Manually designed heuristics are designed to tackle the special requirements such as a bin-item incomplete scenario [14] and multi-tier applications [22, 24]. Although these greedy-based heuristics can quickly approximate an answer, as Mann's research [25] showed, server consolidation is a lot more harder than bin-packing problem because of the multi-dimensional, many constraints. Therefore, general bin-packing algorithms do not perform well with many constraints and specific designed heuristics only perform well in very narrow scope.

Evolutionary Computation (EC) is commonly used to solve combinatorial optimization problem [15], therefore, it is particular useful in solving the global consolidation problem. Many EC techniques including Genetic Algorithm (GA) [50], Ant Colony Optimization (ACO) [13, 27], Particle Swarm Optimization (PSO) [20] have been used in solving this problem. EC algorithms show their advantages in the following aspects. Firstly, EC algorithms are good at solving multi-objective problems because of their population-based nature. And global consolidation problem often involves two or more objectives (e.g energy efficiency and migration cost). Secondly, they can provide near-optimal solutions within a reasonable amount of time.

In VM-based dynamic problem, previous most research proposed human designed greedy-



based dispatching rules or heuristics such as a First-Fit-based approach [5], Modified Best Fit Decreasing [4], and a two-stage heuristic [54]. One of the major problem for human designed heuristics is that if any inherent component gets changes, then the designed heuristic may not work as it was expected [40]. EC algorithms are also seldom considered in this scenario because most EC methods need more time to search through solutions space.

Only a few research focus on container-based consolidation, Piraghaj [33] designs a dynamic allocation system. She proposes a two-step procedure; it first maps tasks to VMs and then allocate containers to VMs. As Mann illustrated in [26], these two steps should be conducted simultaneously, otherwise it leads to local optimal. Other research [11, 17, 1] propose greedy-based heuristics on container allocation problem. They can be easily stuck at local optimal. This thesis, therefore, aims at providing an end-to-end solution for Container-based server consolidation which includes three stages correspond with the Cloud resource management procedure (see Figure 1.6): initialization, static container-based server consolidation and dynamic container placement.

## 1.2 Motivation

The container-based consolidation problem, similar to VM-based consolidation, can be seen as a continuous optimization procedure with three stages: initialization, global joint allocation of container and VM, and container-based dynamic placement. Different stages have distinctive goals, therefore, they are considered as separated research questions. In this thesis, we aims at providing an end-to-end solution to all three problems. In addition, a scalability problem of static optimization is considered as an optional objective.

1. The initialization,  
at this stage, a set of applications or containers is allocated to a set of empty VMs and these VMs are allocated to a set of PMs. This two-step procedure is interconnected meaning the results of first level of allocation affects the second level of allocation. Therefore, both allocations should be conducted simultaneously. This problem is inherently more difficult than previous VM-based consolidation problem. VM-based consolidation is modeled as bin-packing which is NP-hard. In contrast, container-based consolidation has two levels of bin-packing, this is derived from the problem's nature. Most importantly, these two levels of problem interact and therefore can not be solved separately. This is the first research that consider server consolidation has a bi-level optimization problem [48].

This stage will establish the fundamental concepts in studying the joint allocation of containers and VMs including new problem models: price and power model, new problem constraints, and optimization objectives. The major challenges for this objective is to design representations and several EC approaches to solve this problem. More specifically, in designing the EC approach, new search mechanisms, operators will be designed and new representations will be proposed to fit the problem.

2. A *Global consolidation* is conducted to improve the global energy efficiency at a certain time point, e.g. a fixed time interval. The challenges are three folds, firstly, similar with initialization problem, the problem has two level of allocations and they interact with each other. It is more complex than a single-level VM-based consolidation. Secondly, like VM-based consolidation, Container-based consolidation is considered as a multi-objective problem with minimization of migration cost as well as keeping a good energy efficiency. Thirdly, consolidation is a continuous process which means the

previous solution affects the next one. Previous research only consider each consolidation as an independent process. As a consequence, although in current consolidation, the migration is minimized. It may lead to more migration in the future. We will consider the robustness of consolidation and propose a novel time-series-aware server consolidation which takes the previous consolidations and the future consolidation into consideration.

3. Dynamic consolidation continuously maintains the data center to a high energy efficiency. It is applied on single container at any time point. As mentioned in previous Chapter, dynamic placement is directly related to consolidation. Therefore, we focus on this question. To solve a dynamic placement with large number of variables, heuristics and dispatching rules are often used [36, 37, 12, 4]. In this scenario, a dispatching rule is considered as a function that determines the priorities of PMs that a container can be placed. However, dynamic placement is much complex than bin-packing problem [25]. Because of its dynamic nature, human designed heuristics are ill-equipped in approximating solutions when the environment has changed [40]. Multi-objective genetic algorithm (GA) [50] has been applied. However, GA is too slow for dynamic problem.

We intend to develop a hyper-heuristic method - Genetic Programming (GP) technique [2] or artificial immune system [18]- to learn from the best previous allocation and automatic evolves dispatching rules to solve this problem. GP has been applied in generating dispatching rules for bin-packing problem [6, 40] and other scheduling problems [31]. The results have shown promising results.

There are mainly two challenges, first, it is difficult to identify the related factors that construct the heuristic. Factors or features are the building blocks of heuristics. It is a difficult task because the relationship between a good heuristic and features are not obvious. Second, representations provide different patterns to construct dispatching rules. It is also unclear what representation is the most suitable for the consolidation problem.

4. Cloud data center typically has hundreds of thousands PMs and more. Large scale of static server consolidation has always been a challenge since it takes large amount of variables into consider. Many approaches have been proposed in the literature to resolve the problem. There are mainly two ways, both rely on distributed methods, hierarchical-based [21, 30] and agent-based management systems [51]. The major problem in agent-based systems is that agents rely on heavy communication to maintain a high-level utilization. Therefore, it causes heavy load in the networking. Hierarchical-based approaches are the predominate methods. In essence, these approaches are centralized methods where all the states of machines within its region are collected and analyzed. The major disadvantage of hierarchical-based approaches is that it only provides local solutions. In fact, it is infeasible and unnecessary to check all the states of machines since the search space is too large and most machines do not need a change. This idea motivates a way to improving the effectiveness is to reduce the number of variables so that the search space is narrowed. In this thesis, we are going to investigate the way to eliminate the redundant information.

### 1.3 Research Goals

The overall goal of this thesis is to propose an end-to-end server consolidation approach that considers all three stages: Initialization, Off-line Static Joint Allocation of Container and VM,

On-line Dynamic Container Placement Problem. In addition, the static allocation normally involves with large amount of variables which is particular difficult to optimize. We also going to propose a method to solve this problem. These approaches combine element of AI planning, to ensure the objectives and constraint fulfillment, and of Evolutionary Computation, to evolve a population of near-optimal solutions. The research aims to determine a flexible way in creation of solutions to solve server consolidation problems. As discussed in the previous section, the research goal can be achieved in the following objectives and sub-objectives.

1. The initialization Problem,

Currently, most research focus on VM-based server consolidation technique. They often modeled this problem as a vector bin-packing problem [52]. Container adds an extra layer of abstraction on top of VM. The placement problem has become a two-step procedure, in the first step, containers are packed into VMs and then VMs are consolidated into physical machines. These two steps are inter-related to each other.

(a) *Modeling*

Previous VM-based models do not consider two-level allocation structure, therefore, our first sub objective is to propose a description of model for the initialization problem. In order to achieve this goal, we will first review the related models including VM-based placement models and bi-level optimization models. Furthermore, we are going to consider the differences and design the constraints and other characteristics.

(b) *Representation*

Based on this new model, we are going to develop a representation that is suitable for this problem.

(c) *New operators and searching mechanisms*

In order to utilize Evolutionary Computation (EC) to solve this problem, we are going to develop searching mechanisms according to the nature of problem. In order to achieve this goal, we will design several new operators. In order to evaluate the quality of these components, we will perform analytical analysis on the result.

2. Off-line Static Joint Allocation of Container and VM Problem,

A static allocation can be seen as a resource scheduling problem. A schedule is robust if it is able to endure some degree of uncertainty while maintaining a stable solution [8]. Cloud resource management is a continuous process, after each static allocation, the system should be able to maintain a stable status with the least adjustment. The development of static allocation approach has three sub-objectives. In order to measure the degree of robustness, we need to design a robustness measure. The second objective is to design static consolidation algorithm with considering its previous result. The third objective extend the second objective to a more general case, considering both previous and next allocation. The evaluation of algorithm is based on analytical analysis of fitness functions and robustness measure.

(a) *Design a robustness measure*

Previous studies only use simple measurement which counts the migration number between two static consolidation. This measurement aims at minimizing the number of migration in a static placement process. It may cause more migration in the next consolidation. Therefore, it needs a time-series aware measure of the

robustness of system. A data center should be both consolidated as well as robustness after consolidate. Therefore, in this objective, the first sub-problem we are going to solve is to propose a robustness measure.

(b) *Design an allocation method consider previous allocation*

Based on the robustness measure, we will first design an allocation method which takes previous allocation into account. It has two optimization objectives, maximize the robustness and also minimize the energy consolidation.

(c) *Design a time-series-aware allocation method*

Last but not the least, we will generalize the previous sub-objective to a more general one: design a time-series-aware allocation method which takes several allocation into consider.

3. On-line Dynamic Container Placement Problem with a GP approach,

(a) Construct Functional Set and Primitive Set for the problem

As the basic component of a dispatching rule, primitive set contains the states of environment including: status of VMs, features of workloads. The functional set contains the operators which combines low level features.

(b) Representation

In order to utilize a hyper-heuristic method such as GP to solve the problem, the first step is to design a representation.

(c) Develop GP-based methods for evolving Dispatching rules

4. Large-scale Static Consolidation Problem

(a) Propose a preprocessing method to eliminate variables

Current static consolidation takes all servers into consider which will lead to a scalability problem. In this objective, we will propose a method that categorizes servers so that only a small number of servers are considered. This approach will dramatically reduce the search space. The potential approaches that can be applied in this task are various clustering methods.

## 1.4 Published Papers

During the initial stage of this research, some investigation was carried out on the model of container-based server consolidation.

1. Tan, B., Ma, H., Mei, Y. and Zhang, M., "A NSGA-II-based Approach for Web Service Resource Allocation On Cloud". *Proceedings of 2017 IEEE Congress on Evolutionary Computation (CEC2017)*. Donostia, Spain. 5-8 June, 2017.pp.

## 1.5 Organisation of Proposal

The remainder of the proposal is organised as follows: Chapter ?? provides a fundamental definition of the Container-based server consolidation problem and performs a literature review covering a range of works in this field; Chapter ?? discusses the preliminary work carried out to explore the techniques and EC-based techniques for the initialization problem; Chapter ?? presents a plan detailing this projects intended contributions, a project timeline, and a thesis outline.

iiiiii HEAD



## Chapter 2

# Literature Review

### 2.1 Background

This chapter begins by providing an overall understanding of Cloud computing and its related research field. Then, it narrows down to the Energy-aware resource management problem in Section ??.

#### 2.1.1 An Overview of Cloud Computing

Each of the Cloud stakeholders has their goal. End users consume application. They want applications with better quality and lower price. Cloud users develop and deploy softwares on Cloud. They also desire a high Quality of service (QoS) by specifying a Service Level Agreement (SLA) with the Cloud provider.

They are bonded by the Quality of service and computing resource. Quality of service and the computing resources are the two sides of a coin.

For End users, they consume applications deployed on Cloud.

End users consume the application deployed on Cloud. They require a guarantee quality of softwares including functional requirements which are the functionalities defined by Cloud users, and non-functional requirements such as availability, security, and network latency.

Cloud users deploy their software on Clouds. They want to increase the profit by increasing income and decreasing expense. In order to accomplish this goal, they can attract more End users by improving the functionality of softwares and the non-functionality features by guaranteeing Quality of Service (QoS). To improve the non-functionality features, Cloud users need to reserve enough resources as well as minimizing the resource so that the cost is low.

Cloud users identify a list of critical QoS parameters called Service Level Agreement (SLA) which specifies the non-functional requirements such as throughput, latency, and availability. These QoS parameters are mapped to resources (e.g. CPU, memory, network bandwidth) which can satisfy these requirements. Violation of SLA will lead to penalty and decreasing in number of users. Therefore, in essence, the key to attract more users is an effective resource management system which can rapidly react to the fluctuating resource demand.

Beside increase the income, reduce the expense is another way to improve profit. As previous section mentioned, energy consumption is the main source of expense. In energy consumption, server energy consumption is the core that needs to be improved.

### 2.1.2 Energy-aware Resource Management

### 2.1.3 An Overview of Evolutionary Computation

In order to understand Cloud computing, firstly we will illustrate the five essential elements of Cloud computing and their advantages.

Cloud computing has five essential elements:

1. On-demand self-service, it means a Cloud user can require computing resources (e.g CPU time, storage, software use) without the interaction with Cloud provider.
2. Broad network access, Computing resources are connected and delivered over the network.
3. Resource pool, a Cloud provide has a “pool” of resources which are normally virtualized servers. In IaaS, it provides predefined sizes of VMs. In PaaS, the resources are ‘invisible’ to Cloud users who have no knowledge or ability to control.
4. Rapid elasticity, from the perspective of Cloud users, computing resources are assigned and released in real time. In addition, the resources assign to their software is “infinite”. Therefore, Cloud users do not need to worried about the scalability of their applications.
5. Measured Service provides an accurate measure of the usage of computing resources. It is fundamental to the pay-as-you-go policy.

What is your purpose to describe the following content? I would like to discuss the differences, advantages of disadvantage of the resource management in different service models. Therefore, after illustrate how they are work. The point is to compare the resource management. And then, lead to a new service model. And the advantage of new service model should be obvious.

Traditional Cloud computing has three service model as illustrated in Figure ??.

Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service.

### 2.1.4 Resource Management

Scope of Cloud computing resource management.

1. Actors
2. Management Objectives
3. Resource Types
4. Enabling Technologies

Energy-aware Resource Management

This chapter begins by providing a fundamental background to the field of Cloud resource management in Section ??, then addresses several areas of current research interest.



===== 2011cb91584b90ffd83cada4af7a86f40fe3836a



# Bibliography

- [1] ANSELM, J., AMALDI, E., AND CREMONESI, P. Service Consolidation with End-to-End Response Time Constraints. *EUROMICRO-SEAA* (2008), 345–352.
- [2] BANZHAF, W., NORDIN, P., KELLER, R. E., AND FRANCONI, F. D. Genetic programming: an introduction, 1998.
- [3] BARROSO, L. A., AND HÖLZLE, U. The Case for Energy-Proportional Computing. *IEEE Computer* 40, 12 (2007), 33–37.
- [4] BELOGLAZOV, A., ABAWAJY, J. H., AND BUYYA, R. Energy-aware resource allocation heuristics for efficient management of data centers for Cloud computing. *Future Generation Comp. Syst.* 28, 5 (2012), 755–768.
- [5] BOBROFF, N., KOCHUT, A., AND BEATY, K. A. Dynamic Placement of Virtual Machines for Managing SLA Violations. *Integrated Network Management* (2007), 119–128.
- [6] BURKE, E. K., HYDE, M. R., AND KENDALL, G. Evolving Bin Packing Heuristics with Genetic Programming. *PPSN 4193*, Chapter 87 (2006), 860–869.
- [7] 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* 25, 6 (June 2009), 599–616.
- [8] CANON, L.-C., AND JEANNOT, E. Evaluation and Optimization of the Robustness of DAG Schedules in Heterogeneous Environments. *IEEE Transactions on Parallel and Distributed Systems* 21, 4 (2010), 532–546.
- [9] CHAISIRI, S., LEE, B.-S., AND NIYATO, D. Optimization of Resource Provisioning Cost in Cloud Computing. *IEEE Trans. Services Computing* 5, 2 (2012), 164–177.
- [10] COFFMAN JR, E. G., GAREY, M. R., AND JOHNSON, D. S. *Approximation algorithms for bin packing: a survey*. PWS Publishing Co., Aug. 1996.
- [11] DONG, Z., ZHUANG, W., AND ROJAS-CESSA, R. Energy-aware scheduling schemes for cloud data centers on Google trace data. *OnlineGreenComm* (2014), 1–6.
- [12] FORSMAN, M., GLAD, A., LUNDBERG, L., AND ILIE, D. Algorithms for automated live migration of virtual machines. *Journal of Systems and Software* 101 (2015), 110–126.
- [13] GAO, Y., GUAN, H., QI, Z., HOU, Y., AND LIU, L. A multi-objective ant colony system algorithm for virtual machine placement in cloud computing. *J. Comput. Syst. Sci.* 79, 8 (2013), 1230–1242.
- [14] GUPTA, R., BOSE, S. K., SUNDARRAJAN, S., CHEBIYAM, M., AND CHAKRABARTI, A. A Two Stage Heuristic Algorithm for Solving the Server Consolidation Problem with Item-Item and Bin-Item Incompatibility Constraints. *IEEE SCC* (2008).

- [15] GUZEK, M., BOUVRY, P., AND TALBI, E.-G. A Survey of Evolutionary Computation for Resource Management of Processing in Cloud Computing [Review Article]. *IEEE Computational Intelligence ...* 10, 2 (2015), 53–67.
- [16] HAMEED, A., KHOSHKBARFOROUSHHA, A., RANJAN, R., JAYARAMAN, P. P., KOŁODZIEJ, J., BALAJI, P., ZEADALLY, S., MALLUHI, Q. M., TZIRITAS, N., VISHNU, A., KHAN, S. U., AND ZOMAYA, A. Y. A survey and taxonomy on energy efficient resource allocation techniques for cloud computing systems. *Computing* 98, 7 (2016), 751–774.
- [17] HINDMAN, B., KONWINSKI, A., ZAHARIA, M., GHODSI, A., JOSEPH, A. D., KATZ, R. H., SHENKER, S., AND STOICA, I. Mesos - A Platform for Fine-Grained Resource Sharing in the Data Center. *NSDI* (2011).
- [18] HOFMEYR, S. A., AND FORREST, S. Architecture for an Artificial Immune System. *Evolutionary Computation* 8, 4 (2000), 443–473.
- [19] JENNINGS, B., AND STADLER, R. Resource Management in Clouds: Survey and Research Challenges. *Journal of Network and Systems Management* 23, 3 (2015), 567–619.
- [20] JEYARANI, R., NAGAVENI, N., AND RAM, R. V. Design and implementation of adaptive power-aware virtual machine provisioner (APA-VMP) using swarm intelligence. *Future Generation Comp. Syst.* 28, 5 (2012), 811–821.
- [21] JUNG, G., HILTUNEN, M. A., JOSHI, K. R., SCHLICHTING, R. D., AND PU, C. Mistral - Dynamically Managing Power, Performance, and Adaptation Cost in Cloud Infrastructures. *ICDCS* (2010), 62–73.
- [22] JUNG, G., JOSHI, K. R., HILTUNEN, M. A., SCHLICHTING, R. D., AND PU, C. Generating Adaptation Policies for Multi-tier Applications in Consolidated Server Environments. *ICAC* (2008).
- [23] KAPLAN, J. M., FORREST, W., AND KINDLER, N. Revolutionizing data center energy efficiency, 2008.
- [24] LI, B., AND JIANXIN, L. *An Energy-saving Application Live Placement Approach for Cloud Computing Environment*. ACM International Conference on Cloud Computing, 2009.
- [25] MANN, Z. Á. Approximability of virtual machine allocation: much harder than bin packing.
- [26] MANN, Z. Á. Interplay of Virtual Machine Selection and Virtual Machine Placement. In *Service-Oriented and Cloud Computing*. Springer International Publishing, Cham, Aug. 2016, pp. 137–151.
- [27] MATEOS, C., PACINI, E., AND GARINO, C. G. An ACO-inspired algorithm for minimizing weighted flowtime in cloud-based parameter sweep experiments. *Advances in Engineering Software* 56 (2013), 38–50.
- [28] MELL, P. M., AND GRANCE, T. The NIST definition of cloud computing. Tech. rep., National Institute of Standards and Technology, Gaithersburg, MD, Gaithersburg, MD, 2011.
- [29] MISHRA, M., DAS, A., KULKARNI, P., AND SAHOO, A. Dynamic resource management using virtual machine migrations. *IEEE Communications ...* 50, 9 (2012), 34–40.

- [30] MOENS, H., FAMAHEY, J., LATRÉ, S., DHOEDT, B., AND DE TURCK, F. Design and evaluation of a hierarchical application placement algorithm in large scale clouds. *Integrated Network Management* (2011), 137–144.
- [31] NGUYEN, S., ZHANG, M., JOHNSTON, M., AND TAN, K. C. Automatic Design of Scheduling Policies for Dynamic Multi-objective Job Shop Scheduling via Cooperative Coevolution Genetic Programming. *IEEE Transactions on Evolutionary Computation* 18, 2 (2014), 193–208.
- [32] PANIGRAHY, R., TALWAR, K., UYEDA, L., AND WIEDER, U. Heuristics for vector bin packing. *research.microsoft.com* (2011).
- [33] PIRAGHAJ, S. F., CALHEIROS, R. N., CHAN, J., DASTJERDI, A. V., AND BUYYA, R. Virtual Machine Customization and Task Mapping Architecture for Efficient Allocation of Cloud Data Center Resources. *Comput. J.* 59, 2 (2016), 208–224.
- [34] PIRAGHAJ, S. F., DASTJERDI, A. V., CALHEIROS, R. N., AND BUYYA, R. Efficient Virtual Machine Sizing for Hosting Containers as a Service (SERVICES 2015). *SERVICES* (2015).
- [35] RONG, H., ZHANG, H., XIAO, S., LI, C., AND HU, C. Optimizing energy consumption for data centers. *Renewable and Sustainable Energy Reviews* 58 (May 2016), 674–691.
- [36] SARIN, S. C., VARADARAJAN, A., AND WANG, L. A survey of dispatching rules for operational control in wafer fabrication. *Production Planning and ...* 22, 1 (Jan. 2011), 4–24.
- [37] SHI, W., AND HONG, B. Towards Profitable Virtual Machine Placement in the Data Center. *UCC* (2011), 138–145.
- [38] SOLTESZ, S., PÖTZL, H., FIUCZYNSKI, M. E., BAVIER, A. C., AND PETERSON, L. L. Container-based operating system virtualization - a scalable, high-performance alternative to hypervisors. *EuroSys* 41, 3 (2007), 275–287.
- [39] SOMANI, G., AND CHAUDHARY, S. Application Performance Isolation in Virtualization. *IEEE CLOUD* (2009), 41–48.
- [40] SOTELO-FIGUEROA, M. A., SOBERANES, H. J. P., CARPIO, J. M., HUACUJA, H. J. F., REYES, L. C., AND SORIA-ALCARAZ, J. A. Evolving Bin Packing Heuristic Using Micro-Differential Evolution with Indirect Representation. *Recent Advances on Hybrid Intelligent Systems* 451, Chapter 28 (2013), 349–359.
- [41] SPEITKAMP, B., AND BICHLER, M. A mathematical programming approach for server consolidation problems in virtualized data centers. *IEEE Transactions on services ...* (2010).
- [42] SVARD, P., LI, W., WADBRO, E., TORDSSON, J., AND ELMROTH, E. Continuous Datacenter Consolidation. In *2015 IEEE 7th International Conference on Cloud Computing Technology and Science (CloudCom)* (2015), IEEE, pp. 387–396.
- [43] TOMÁS, L., AND TORDSSON, J. Improving cloud infrastructure utilization through overbooking. *CAC* (2013), 1.
- [44] UHLIG, R., NEIGER, G., RODGERS, D., SANTONI, A. L., MARTINS, F. C. M., ANDERSON, A. V., BENNETT, S. M., KÄGI, A., LEUNG, F. H., AND SMITH, L. Intel Virtualization Technology. *IEEE Computer* 38, 5 (2005), 48–56.

- [45] VERMA, A., AND DASGUPTA, G. Server Workload Analysis for Power Minimization using Consolidation. *USENIX Annual Technical Conference* (2009), 28–28.
- [46] WANG, Y., AND XIA, Y. Energy Optimal VM Placement in the Cloud. In *2016 IEEE 9th International Conference on Cloud Computing (CLOUD)* (2016), IEEE, pp. 84–91.
- [47] WEI, Y., AND BLAKE, M. B. Service-Oriented Computing and Cloud Computing - Challenges and Opportunities. *IEEE Internet Computing* 14, 6 (2010), 72–75.
- [48] WEN, U.-P., AND HSU, S.-T. Linear Bi-Level Programming Problems – A Review. *The Journal of the Operational Research Society* 42, 2 (Feb. 1991), 125.
- [49] XIAO, Z., JIANG, J., ZHU, Y., MING, Z., ZHONG, S.-H., AND CAI, S.-B. A solution of dynamic VMs placement problem for energy consumption optimization based on evolutionary game theory. *Journal of Systems and Software* 101 (2015), 260–272.
- [50] XU, J., AND FORTES, J. A. B. Multi-Objective Virtual Machine Placement in Virtualized Data Center Environments. *GreenCom/CPSCoM* (2010).
- [51] YAZIR, Y. O., MATTHEWS, C., FARAHBOD, R., NEVILLE, S. W., GUITOUNI, A., GANTI, S., AND COADY, Y. Dynamic Resource Allocation in Computing Clouds Using Distributed Multiple Criteria Decision Analysis. *IEEE CLOUD* (2010), 91–98.
- [52] ZHANG, J., HUANG, H., AND WANG, X. Resource provision algorithms in cloud computing - A survey. *J. Network and Computer Applications* 64 (2016), 23–42.
- [53] ZHANG, Q., CHENG, L., AND BOUTABA, R. Cloud computing: state-of-the-art and research challenges. *Journal of Internet Services and Applications* 1, 1 (2010), 7–18.
- [54] ZHANG, Z., HSU, C.-C., AND CHANG, M. Cool Cloud: A Practical Dynamic Virtual Machine Placement Framework for Energy Aware Data Centers. In *2015 IEEE 8th International Conference on Cloud Computing (CLOUD)* (2015), IEEE, pp. 758–765.