

# Chapter 23

## Resource Provisioning in the Cloud: An Exploration of Challenges and Research Trends

**Ming Mao**

*University of Virginia, USA*

**Marty Humphrey**

*University of Virginia, USA*

### **ABSTRACT**

*It is a challenge to provision and allocate resources in the Cloud so as to meet both the performance and cost goals of Cloud users. For a Cloud consumer, the ability to acquire and release resources dynamically and trivially in the Cloud, while being a powerful and useful aspect, complicates the resource provisioning and allocation task in the Cloud. While on the one hand, resource under-provisioning may hurt application performance and deteriorate service quality; on the other hand, resource over-provisioning could cost users more and offset Cloud advantages. Although resource management and job scheduling have been studied extensively in the Grid environments and the Cloud shares many common features with the Grid, the mapping from user objectives to resource provisioning and allocation in the Cloud has many challenges due to the seemingly unlimited resource pools, virtualization, and isolation features provided by the Cloud. This chapter focuses on surveying the research trends in resource provisioning in the Cloud based on several factors such as the type of the workload, the VM heterogeneity, data transfer requirements, solution methods, and optimization goals and constraints, and attempts to provide guidelines for future research.*

DOI: 10.4018/978-1-4666-6178-3.ch023

## INTRODUCTION

The Cloud has become a significant computing platform. It has attracted many businesses and individual users by offering on-demand computing power and storage capacity. The economies of scale and *pay-as-you-go billing model* could save users large up-front capital investments and long term operation costs (Armbrust et al., 2010). A key feature of the Cloud is the *elasticity*, the ability to dynamically acquire and release computing resources in response to demand (Mell & Grance 2011). For successful resource management in the Cloud, one needs to first determine and acquire the required amount and type of resources that may be needed to satisfy a computing job, which is the act of *provisioning*; and then place computing activities onto each of the resources in a dynamic and efficient manner, which is the act of *allocation*. The provisioning and allocation of Cloud resources is a challenging problem because the mapping from user objectives to the resource provisioning and allocation plans is not trivial (Mell & Grance, 2011; Buyya et al., 2009).

Resource provisioning and allocation in the Cloud needs to consider the following factors:

- A performance goal may be achieved through different types of resources with different costs.
- A fixed budget may be used to rent a wide variety of resource configurations for varying durations.
- Task precedence constraints may need to be preserved for a job.
- The workload may experience unexpected peaks.
- The performance requirements and cost constraints may change dynamically.

One of the main advantages of the Cloud is the rapid elasticity or dynamic scalability (Armbrust

et al., 2010), which enables the dynamic acquisition and release of Cloud resources in response to demand. It is a key enabler of Cloud adoption. This elasticity saves the Cloud users large up-front capital investments and allows the computing resources to grow according to business demand. The elasticity has become one of the main forces to drive application migration to the Cloud.

Another important factor when considering Cloud adoption is the cost aspect (Armbrust et al., 2010). Maximizing the return and minimizing the cost of Cloud investment are the two main goals for Cloud users. The cost savings result from the economies of scale and dynamic scalability. For the same business scenarios, if it is more expensive to develop and maintain the applications in the Cloud, the Cloud loses its advantages. Therefore, cost savings is a significant concern for Cloud users. In addition to being a goal, sometimes cost could become a constraint for some Cloud applications. For example, Cloud consumers may have a budget limit that they are allowed to spend on Cloud purchasing which restricts the running cost of the acquired resources from exceeding a certain amount. In such cases, cost essentially determines the maximum size of the acquired resource pool.

In summary, the key benefit of the Cloud is to be able to acquire resources in response to demand dynamically and only pay for the resources used. This benefit can only be realized when the Cloud users can determine the right size of the resource pool and allocate the resources in a cost-effective way. While resource over-provisioning can cost users more than necessary, which essentially offsets the Cloud advantages; resource under-provisioning hurts the application performance and could violate the service level agreements that service providers on the Cloud have with their customers, causing customers to turn away. Essentially, the Cloud adopters should understand what resources should be acquired or released

22 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

[www.igi-global.com/chapter/resource-provisioning-in-the-cloud/115445](http://www.igi-global.com/chapter/resource-provisioning-in-the-cloud/115445)

## Related Content

---

### Improving the Detection of On-Line Vertical Port Scan in IP Traffic

Christine Fricker, Philippe Robertand Yousra Chabchoub (2014). *International Journal of Secure Software Engineering* (pp. 61-74).

[www.irma-international.org/article/improving-the-detection-of-on-line-vertical-port-scan-in-ip-traffic/109581](http://www.irma-international.org/article/improving-the-detection-of-on-line-vertical-port-scan-in-ip-traffic/109581)

### Graphical-Based Authentication System and Its Applications

Priti Golarand Brijesh Khandelwal (2021). *Design, Applications, and Maintenance of Cyber-Physical Systems* (pp. 63-91).

[www.irma-international.org/chapter/graphical-based-authentication-system-and-its-applications/281769](http://www.irma-international.org/chapter/graphical-based-authentication-system-and-its-applications/281769)

### Adaptive Future Internet Applications: Opportunities and Challenges for Adaptive Web Services Technology

Clarissa Cassales Marquezan, Andreas Metzger, Klaus Pohl, Vegard Engen, Michael Boniface, Stephen C. Phillipsand Zlatko Zlatev (2018). *Application Development and Design: Concepts, Methodologies, Tools, and Applications* (pp. 1568-1589).

[www.irma-international.org/chapter/adaptive-future-internet-applications/188271](http://www.irma-international.org/chapter/adaptive-future-internet-applications/188271)

### Determination of Melting Point of Chemical Substances Using Image Differencing Method

Anurag Shrivastavaand Rama Sushil (2022). *International Journal of Software Innovation* (pp. 1-10).

[www.irma-international.org/article/determination-of-melting-point-of-chemical-substances-using-image-differencing-method/297985](http://www.irma-international.org/article/determination-of-melting-point-of-chemical-substances-using-image-differencing-method/297985)

### A State-Based Intention Driven Declarative Process Model

Pnina Soffer (2013). *International Journal of Information System Modeling and Design* (pp. 44-64).

[www.irma-international.org/article/state-based-intention-driven-declarative/80244](http://www.irma-international.org/article/state-based-intention-driven-declarative/80244)