# Predicting provisioning and booting times in a Metal-as-a-service system

Alexandru Sîrbu [a], Cristian Pop [a], Cristina Şerbănescu [b], Florin Pop [a,*]

[a] *Computer Science Department, Faculty of Automatic Control and Computers, University Politehnica of Bucharest, Romania*
[b] *Department of Mathematical Methods and Models, Faculty of Applied Sciences, University Politehnica of Bucharest, Romania*

## HIGHLIGHTS

- Architecture, entities and operations for a MaaS system.
- Provisioning flow and stages.
- Regression-based prediction algorithm for booting time.
- Performance evaluation of instances and cluster provisioning.

## ARTICLE INFO

## ABSTRACT

Cloud management automation and management of SLA incidents become a research challenges for any Cloud service-based system. In the era of ongoing adoption of Cloud Computing at a fast rate the Metal-as-a-service (MaaS) platforms assure a higher level of performance, but at the cost of a more complex provisioning system, all of these being imposed by SLA assurance. More, disaster recovery and critical infrastructure protection become important aspects for any real-time applications that use Cloud Services. This paper deals with the problem of predicting provisioning and booting times in a MaaS system, and proposed a solution based on platform monitoring and a multi-variate regression algorithm. The configuration, provisioning flow, and capacity management capabilities were tested on Bigstep Full Metal Cloud platform an event-based tracking system, based on which provisioning times can be calculated for each individual element. We analyzed the performance of proposed solution by comparing the predicted booting and provisioning times with real times using different scenarios.

## 1. Introduction

The basic building block of a Cloud solution consists of virtualization, in which the physical servers are shared by multiple tenants utilizing virtual machines, which access the hardware resources as stated by the Service-Level Agreement (SLA) [1]. In the new fashion of metal-as-a-service (MaaS), which came as a performance improvement over the classic cloud, the end user has full processing power of the physical server, without using a virtualization layer or sharing the resources with other users. Resource clustering [2] in datacenters is one solution used in allocation problems that need to follow a specific SLA. In Clouds, one important aspect is that using an external provisioning system, the time it takes from requesting a server until receiving one is higher than the time a virtual machine is cloned from a template and given to a client. This aspect is very important for Big Data processing platform that need to support a large number of tasks [3].

In a time in which SLAs manage the relationship between the cloud provider and client, knowing the amount of time required for starting a new instance (taking into consideration the actual provisioning time and, in the case of a bare metal cloud, the booting time of the server, which is not ignorable any more as opposed to a virtualized cloud) is mandatory for both entities. The client needs to know what to expect from the provider and when their services will be up and running and the provider must take into consideration the times when signing an SLA, in order to actually be able to deliver the promised performances. Sometimes this SLA can be automatically established for cloud computing services using a policy-based system [4]. This is why a prediction system implemented for approximating this duration is one of the important concerns when talking about cloud provisioning and is

* Corresponding author. Fax: +40 318 145 309.
*E-mail address:* florin.pop@cs.pub.ro (F. Pop).

even more important for metal-as-a-service providers, who use more services and actual hardware provisioning [5–7].

In this growing world of bare metal cloud solutions, more and more services are added coming from its virtualized counterpart, and one of these is the usage of a reservation and advanced provisioning system in which clients can request a server for a certain period of time in advance, knowing they need to use it during that period [8]. As provisioning takes longer on a metal infrastructure, a client needs to carefully choose the periods in which he wants to request a server, in order to use its computing power and to minimize his monetary investment [9]. The paper studies this problem, by testing the provisioning capabilities of a metal as a service provisioning system with a series of common usage scenarios (like advanced reservation system), giving some time-gaining and money-saving solutions.

Lately, the concept of hybrid clouds promised to solve the problem of limited elasticity of the private clouds [10]. According to its principles, enterprises may still run their own private clouds for their sensitive tasks, but the less sensitive ones may be executed on public clouds and usage spikes may also be alleviated by spawning additional machines in the public cloud [11]. The bare-metal cloud comes as an alternative to the hybrid clouds by allowing the enterprises to run their own private clouds on the provider's infrastructure without the burden of acquiring and installing additional hardware when needed. The private clouds may be allowed to scale easily when additional processing power is needed. In this case, the bare-metal cloud acts as a cloud of clouds.

Also, when compared with the classical grid/cluster computing paradigm, the fact that the physical resources (servers, network equipment) are allocated and programmed dynamically by a centralized system and without doing any changes to the datacenter topology (thus, allowing them to change their purpose in the system when they need to) is the main addition and difference, allowing grids and other clusters to be provisioned inside a bare metal cloud.

With the implementation of an advanced reservation system which minimizes the time between provision finalization and utilization, the provider can anticipate the duration of many of the common operations which can occur in the system (provisioning, booting, de-provisioning) and can offer new clauses in its SLA, regarding the delivery time of its services. This is important, as metal-as-a-service provisioning times are larger than the ones in classic clouds (as it implies also working on the actual hardware and not only software-based), but are notably lower than ordering dedicated servers, for example (as many parts are automated and networking is taken cared of automatically and allows more flexibility) [12]. Thus, a client can know beforehand the time needed for service operations and can rely on them, in order to maximize its resources (both computing power and financial), and the provider can use the actual servers more efficiently, as the resources on the provider's part have to be used with greater care, as the actual scaling is limited to the hardware infrastructure.

The above aspects were presented in the previous published paper, "MaaS Advanced Provisioning and Reservation System" [13]. The current paper expands those ideas, in order to have a more complex prediction engine, which generates more exact information to clients and service providers, with a practical application in a bare metal cloud, by having the following main contributions:

- we present an extended version of Bigstep Full Metal Cloud[1] solution by describing the main e entities and operations;
- we consider provisioning stages, which highlights the stages involved in the creation of the building block of the bare metal Cloud;

- we used the prediction engine for a bare metal cloud provider, which takes into consideration the provisioning graph for a given deployment and the duration of previous similar tasks;
- we proposed a prediction algorithm for the boot duration that uses a divide and conquer approach to filter the past data and a linear regression to provide estimations, and also an extended model for multiple regression, which is applied on various factors that affect the boot duration for any MaaS platform;
- we evaluate the proposed estimation solution with real-time measurements and we provide the building blocks for a more complex advanced reservation system.

The paper is organized as follows. In Section 2, the advanced reservation and provisioning problem is defined. Then, in Section 3, a series of MaaS solutions are presented as related-work solutions. Section 4 presents in detail the Bigstep Full Metal Cloud and the provisioning flow. Afterwards, the tests and results are presented and then Section 5 describes a possible approach to provisioning duration prediction based on linear regression, alongside a series of tests and results, which could be utilized to implement such a system. The paper ends with the conclusions and some possible future work directions, presented in Section 6.

## 2. Problem definition and related work

### 2.1. Problem definition

Some cloud users require servers with plain operating systems to run their own applications, while others, especially those coming from the Big Data world, use commercial or open source distributed applications like Datastax.[2] In both situations, there are cases when the need or even actual quantity of computing resources can be predicted. In such cases, a reservation and advanced provisioning system that takes the predicted start time and the quantity of resources as input and makes sure the resources are available and provisioned when needed, could prove to be a useful tool.

The availability of resources is a simple problem, as it can be solved by marking them as reserved. Although this approach could prove to be problematic, as it might keep resources locked for large time periods in which these may be needed elsewhere, the availability is not the object of this study and as a result the simplest solution suffices.

The most elementary reservation and advanced provisioning system could start the provisioning of the resources at the time specified by the user, but this would require him to be familiar with the internals of the used provisioning system in order to compensate for the delays it causes and this is usually not the case [5,14]. Therefore, a reservation system should be aware of these delays and their causes in order to reduce the idle time of the provisioned resources. Before a server can be booted various provisioning steps that may cause significant delays have to be executed such as creating the external storage devices with pre-installed operating systems, configuring the network to provide isolation and secure communication, populating the DNS servers with appropriate records and others. Afterwards, the server can be powered on and the operating system is booted through the network. The booting duration may increase considerably due to network congestion or storage server load as these may lead to an increased communication latency or even packet loss. As these distinct phases of the provisioning process are causing non-constant delays that vary depending on different variables, an

---

[1] http://bigstep.com/full-metal-cloud.

[2] www.datastax.com.