



A meta-heuristic optimization approach to the scheduling of bag-of-tasks applications on heterogeneous clouds with multi-level arrivals and critical jobs



Ioannis A. Moschakis*, Helen D. Karatza

Department of Informatics, Aristotle University of Thessaloniki, 54124 Thessaloniki, Greece

ARTICLE INFO

Article history:

Received 17 February 2015

Received in revised form 6 April 2015

Accepted 7 April 2015

Available online 8 June 2015

Keywords:

Tabu-search

Simulated annealing

Bag-of-tasks

Multi-criteria scheduling

ABSTRACT

As cloud computing evolves, it is becoming more and more apparent that the future of this industry lies in interconnected cloud systems where resources will be provided by multiple “Cloud” providers instead of just one. In this way, the hosts of services that are cloud-based will have access to even larger resource pools while at the same time increasing their scalability and availability by diversifying both their computing resources and the geographical locations where those resources operate from. Furthermore the increased competition between the cloud providers in conjunction with the commoditization of hardware has already led to large decreases in the cost of cloud computing and this trend is bound to continue in the future. Scientific focus in cloud computing is also headed this way with more studies on the efficient allocation of resources and effective distribution of computing tasks between those resources. This study evaluates the use of meta-heuristic optimization algorithms in the scheduling of bag-of-tasks applications in a heterogeneous cloud of clouds. The study of both local and globally arriving jobs has been considered along with the introduction of sporadically arriving critical jobs. Simulation results show that the use of these meta-heuristics can provide significant benefits in costs and performance.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Cloud computing is no longer in its infancy, it has evolved from a niche of early adopters and web-startups into a robust framework of technologies and has established itself as a major computing platform. It has also managed to attain a large share of the distributed computing market with more and more companies using the cloud to offer their services or process their data at a fraction of the cost that maintaining an in-house IT infrastructure would have costed in the past, effectively lowering the barrier of entry to such businesses. Moreover, its ability to scale “on-demand” while providing a flexible cost model coupled with high availability and ease of customization, positions it ahead of classic distributed technologies such as the Grid [1,2].

Clouds are also multifaceted, if one requires just computing capacity and wants to customize it according to his needs the infrastructure clouds can offer that ability in the form of virtual machines (VMs) or, more recently, containers. Others may just need a coherent framework on which to build their work without the need to care about lower-level details required by infrastructure clouds, Platform-as-a-Service and Software-as-a-Service clouds are able to fill in these requirements. Finally, a

* Corresponding author.

E-mail addresses: imoschak@csd.auth.gr (I.A. Moschakis), karatza@csd.auth.gr (H.D. Karatza).

company may need to outsource parts of its computing capacity while also maintaining other parts in-house. In this case Private and Hybrid clouds can provide a solution [3].

One of the major challenges, identified even from the beginning of the cloud computing industry, is the inter-linking of clouds, where resources from multiple clouds provided by different providers and possibly different architectures are coalesced transparently into a single entity. Recently, there has been much scientific effort into developing models allowing for such combinations of heterogeneous resources not only from a performance perspective but also from an architectural standpoint.

This kind of conglomeration presents many obstacles involving the use of open communication standards and middle-ware frameworks, which act as glue between the different cloud systems. Also, traditional distributed schedulers may prove to be insufficient for the scheduling of such systems which consist of highly scalable, geographically distributed heterogeneous systems and therefore warrant the introduction of more modern schedulers [4].

When multiple cloud providers are mixed together, traditional single layered scheduling approaches typical in distributed systems will not suffice [5]. In order to control the use of multiple clouds a multi-layered approach is required using a global scheduler that controls the distribution of jobs between clouds as well as local to the cloud schedulers which determine how each cloud's resources should be allocated to incoming jobs [6].

The scheduler's must also be aware of several factors in order to make informed decisions on how to partition workload between resources. Merging VMs of several clouds is bound to hit heterogeneity issues which will affect both the performance and the capacity of the cloud. Cost is also another important factor which will come into play and may need to be continuously re-examined, as resources are leased and released, as cloud costs are not constant and usually tend to change over time. In addition, individual clouds may also service their private customers which may be given higher priority than those coming from the multi-cloud. It is important therefore to examine the effect that these "local" jobs may have on the multi-cloud system.

The nature of the jobs being serviced is an equally important factor of such a model. "Big-data" applications use the cloud for processing large volumes of data in parallel and recombine results once the individual parts have finished. The MapReduce processing model [7] is arguably the biggest example of such a method of processing data. A very well known and studied model of parallel jobs that fits these requirements is the Bag-of-Tasks (BoT) model [8]. In this model of embarrassingly parallel jobs, individual job tasks are allowed to execute in any order and in any available resource. The completion of the last task marks also the completion of the entire BoT.

The model presented in this study, considers the application of two optimization meta-heuristic methods, namely Simulated Annealing (SA) and Tabu-Search (TS) when applied to the scheduling of a heterogeneous multi-cloud system with multi-layered arrivals and critical jobs. The choice of using meta-heuristic methods, and comparing them with currently available state-of-the-art heuristic methods for the same problem, was driven by the complexity of the model in study as well as the computational hardness of BoT scheduling systems in general. As we will see later in this paper, and judging from the cost and performance results, this choice is well justified.

For the generation of the BoT synthetic workload traces, used in the simulation experiments for this model, we have used the framework introduced by Minh et al. in [9–11]. This framework was developed with BoT specific traits and characteristics in mind, such as long range dependence (LRD), job periodicity and temporal burstiness. Therefore, the synthetic traces produced by this framework are representative of and aligned with the workload model considered this study.

The rest of this paper is structured as follows. In Section 2 we provide references to relative works and compare them with ours. In Section 3 we analyze the cloud system model that was implemented. Section 4 describes the model of BoT applications. In Section 5 we discuss the dispatching and scheduling processes used by the system along with the algorithms applied. Section 6 presents the metrics used to measure the performance and cost of the system. Finally in Sections 7 and 8 we present the results along with an analysis of them and provide conclusive remarks and thoughts about our future work.

2. Related work

Many recent models have applied the use of bag-of-tasks applications in clouds. In [12] Oprescu et al. employ a tail-phase optimization algorithm for the minimization of cost taking advantage of idling holes in the scheduling window to maximize resource utilization and minimize costs at the same time. In [13] the same authors present their "BaTS" scheduler which is used to calculate the cost and performance of a BoT on multiple different clouds under budget-constraints and offer viable options to the user. Farahabady et al. in [14] introduce their FPRAS algorithm in the assignment of BoT applications in multiple clouds. Netto and Buyya in [15] propose multiple policies for evaluating resource offerings from cloud providers to schedule deadline-constrained BoT applications, while in [16] they examine a coordinated rescheduling strategy used to deal with imprecise run-time estimates for BoTs running on multiple clouds. Finally Garcia and Sim in [17] employ the use of a genetic algorithm for estimating resource sets required for a BoT execution coupled with an agent-based system for executing time and budget constrained BoTs.

Apart from the cloud, the use of bag-of-tasks scheduling has been extensively studied in the context Grid [8,18,19] and Distributed Systems [20,21]. Various models have been developed incorporating features like power efficiency [22], deadlines [23] and fault tolerance [24,25]. However none of these studies considered the cost as an important factor since it is not a traditional characteristic of Grid and Distributed systems. Also the systems examined were static without any of the dynamic characteristics found in cloud models.

Download English Version:

<https://daneshyari.com/en/article/491911>

Download Persian Version:

<https://daneshyari.com/article/491911>

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