

Multi-Resource Workload Consolidation in Cloud Data Centers

Carlo Mastroianni
ICAR-CNR and
eco4cloud srl
via P Bucci 41C
87036 Rende (CS), Italy
Email: mastroianni@icar.cnr.it

Michela Meo
Department of Electronics and
Communications, Politecnico di Torino
Corso Duca degli Abruzzi, 24
10129 Torino, Italy
Email: michela.meo@polito.it

Giuseppe Papuzzo
ICAR-CNR and
eco4cloud srl
via P Bucci 41C
87036 Rende (CS), Italy
Email: papuzzo@eco4cloud.com

Abstract—Consolidation of Virtual Machines (VMs) on the minimum number of physical servers has been recognized as a very efficient approach to increase the efficiency of virtualized data centers and save energy, as consolidation allows unloaded servers to be switched off or used to accommodate more load. The problem is so complex that centralized and deterministic solutions are useless in large data centers with hundreds or thousands of servers. This paper presents a self-organizing approach for the consolidation of VMs on two resources, CPU and RAM. Decisions on the assignment and migration of VMs are driven by probabilistic processes and are based on local information, which makes the solution simple to implement and scalable. Experiments on a real data center show that the approach rapidly consolidates the workload, and CPU-bound and RAM-bound VMs are balanced, so that both resources are exploited efficiently.

I. INTRODUCTION

The ever increasing demand for computing resources has led companies and Cloud providers to build and use large data centers, which consume a significant amount of energy (about 1.5% of the produced electricity) and cause huge carbon emissions [1]. Consolidation of the workload, which consists in allocating the maximum number of VMs in the minimum number of physical machines [2], helps to alleviate the problem. After consolidation, unneeded servers can be switched off or devoted to the execution of incremental workload. Unfortunately, efficient VM consolidation is hindered by the inherent complexity of the problem. The optimal assignment of VMs to the servers of a data center is analogous to the NP-hard “Bin Packing Problem”, the problem of assigning a given set of items of variable size to the minimum number of bins taken from a given set.

In [3] we presented ecoCloud, an approach for consolidating VMs according to a single computing resource, i.e., the CPU. The approach has been extended to the multi-dimension problem, and specifically for the case in which VMs are consolidated with respect to two resources: CPU and RAM.

With ecoCloud, VMs are consolidated using two types of probabilistic procedures, for the *assignment* and the *migration* of VMs. Both procedures aim at increasing the utilization of servers and consolidating the workload dynamically, with the twofold objective of saving electrical costs and respecting the contracts stipulated with users, especially concerning the

expected quality of service. The scenario is pictured in Figure 1: the request to run a client application is transmitted to the data center manager, which selects an appropriate VM and assigns the VM to one of the available servers through the *assignment procedure*.

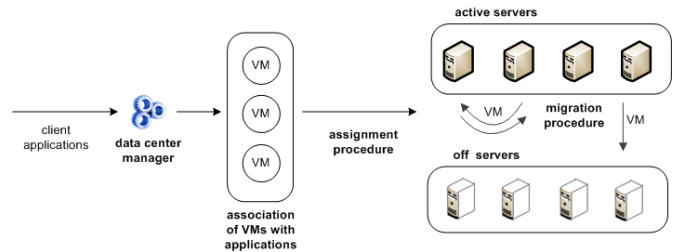


Fig. 1. Assignment and migration of VMs in a data center.

Upon an invitation from the central manager, a single server autonomously decides whether to give or deny its availability to accept a VM. Decisions are based on local information about the utilization of CPU and RAM, and are founded on Bernoulli trials, whose success probability is driven by a probabilistic function. The function is defined so as to favor the acceptance of a VM in servers with intermediate utilization, as this helps to foster consolidation, while the invitation tends to be rejected by over-utilized and under-utilized servers. In the case of over-utilization, the rationale is to avoid overload situations that can penalize the quality of service perceived by users, while in the case of under-utilization the objective is to try to get rid of the running VMs and switch off the server. The data center manager assigns the VM to one of the servers that have given their availability.

The workload of VMs changes with time: therefore, the assignment of VMs is monitored continuously and is tuned through the *migration procedure*, also driven by Bernoulli trials. A migration of a VM to another server is requested with high probability either when the utilization of server resources is too low, meaning that the server is under-utilized, or when it is too high. Details on the assignment and migration procedures are given in [4]. The following section describes the results obtained in a real data center.

II. EXPERIMENTS ON A REAL DATA CENTER

This section reports the results of the experiments performed in July 2013 on a data center owned by an Italian telecommunication company. The experiment was run on 28 servers virtualized with the platform VMWare vSphere. The 447 VMs hosted by the data center were categorized into CPU-bound (C-type) and memory-bound (M-type) depending on their usage of the two resources. In this data center, 75 percent of the VMs, 335, were memory-bound and 112 were CPU-bound. The M-type VMs contributed for about 40% of the overall CPU load and 90% of the overall memory load.

Figure 2 shows the number of active servers starting from the time at which ecoCloud is activated and for the following 12 days. Within the first three days 11 servers, out of 28, are switched off. Figure 3 reports the number of migrations from over- and under-utilized servers (referred to as “high” and “low” migrations, respectively) performed during the analyzed period on the whole data center. In the first days, migrations are mostly from low utilized servers, which are first unloaded and then switched off. As the consolidation process proceeds, the number of migrations stabilizes to definitely acceptable values: for example, in the last two days no more than four migrations per day are performed.

Figure 4 offers a snapshot of the data center at the end of the twelfth day of ecoCloud operation, when only 17 of 28 are still active. The figure reports, for each of the 28 servers, the amount of CPU and RAM utilized by C-type and M-type VMs. Since in this scenario most VMs are memory-bound, the consolidation is driven by RAM: in the majority of active servers the RAM utilization is over 70%, three servers have a RAM utilization between 60% and 70%, and a single server – the one labeled as server 6 – has a RAM utilization lower than 50%. The consolidation is made possible by the fact that VMs of the two types are distributed among the servers in a proportion that never diverts too much from the overall proportion observed in the whole data center.

III. CONCLUSION

The paper focuses on the problem of making data centers and Cloud infrastructures more energy efficient. One of the

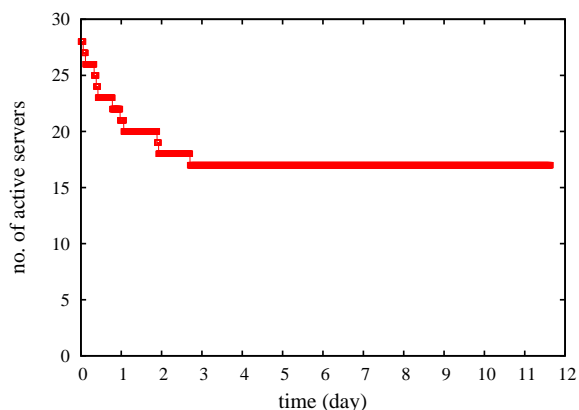


Fig. 2. Number of active servers after activation of ecoCloud.

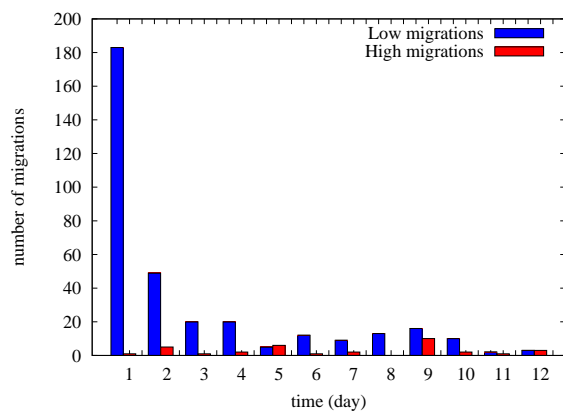


Fig. 3. Number of VM migrations after activation of ecoCloud.

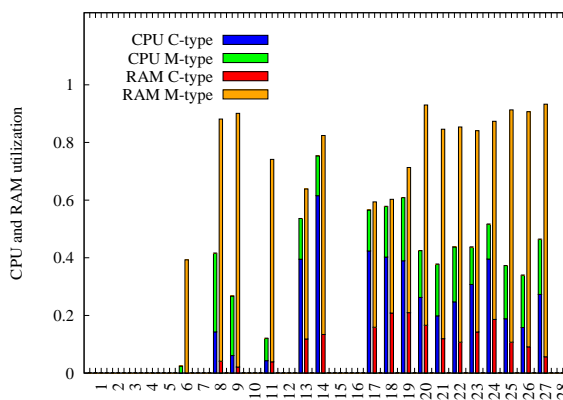


Fig. 4. RAM and CPU utilization on the 28 servers, separated for the C-type and M-type VMs. Values are taken at the end of the 12th day of operation.

most promising approaches consists in dynamically consolidating the load on as few servers as possible. In particular, the paper deals with a recently proposed solution, namely ecoCloud, which, by being decentralized and probabilistic in nature, is highly scalable and allows smooth adaptation of the infrastructure to the actual traffic load. Experiments on a real data center show that ecoCloud achieves high consolidation and effectively balances CPU-bound and RAM-bound VMs. More details and experiments, as well as a thorough mathematical analysis, can be found in [4].

REFERENCES

- [1] A. Beloglazov, J. Abawajy, and R. Buyya, “Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing,” *Future Generation Computer Systems*, vol. 28, no. 5, pp. 755–768, 2012.
- [2] P. Graubner, M. Schmidt, and B. Freisleben, “Energy-efficient virtual machine consolidation,” *IT Professional*, vol. 15, no. 2, pp. 28–34, 2013.
- [3] C. Mastroianni, M. Meo, and G. Papuzzo, “Self-economy in cloud data centers: Statistical assignment and migration of virtual machines,” in *17th International European Conference on Parallel and Distributed Computing, Euro-Par 2011*, vol. 6852. Bordeaux, France: Springer LNCS, September 2011, pp. 407–418.
- [4] —, “Multi-resource workload consolidation in cloud data centers,” ICAR-CNR, Tech. Rep., September 2013. [Online]. Available: <http://www.icar.cnr.it/tr/2013/02>