# An Automated Approach towards Optimal Resource Allocation in Virtualized Data Center

Sajib Kundu, Raju Rangaswami, Kaushik Dutta
Florida International University
skund001@cs.fiu.edu, raju@cs.fiu.edu, kaushik.dutta@business.fiu.edu

## 1   Introduction

Server consolidation is becoming a common practice in the data centers. There are several reasons for doing this. Firstly, it reduces the power consumption which cuts down the recurring costs. Secondly, it decreases the setup expenditures by allowing fewer physical resources to host more number of applications. Thirdly, it also reduces the supervising effort of the administrators to a great extent. Originally, data centers used to allow each customer-preferred application to have its own dedicated processing unit, memory, network and storage devices. Due to the strong requirements imposed by the Service-Level Agreements (SLAs) and the unavailability of any good consolidation solutions, often times the servers were under-utilized. The rapid growth of virtualization technologies is making the objective of server consolidation feasible and readily achievable. Data centers are now reducing power consumption and hardware costs by running several virtual machines within a single physical host. As virtualization allows strong resource partitioning among virtual machines (VMs), applications running inside those VMs are still guaranteed the desired Quality-of-Service (QoS). As more or more VMs are confined to operate inside the same pieces of hardware, the need for an efficient resource distribution policy is also increasing. There is also an underlying objective of augmenting the utilization of the resources thereby, maximizing the profit for the data centers. The task of allocating resources is often manual which are derived from some high-level policies set up by the administrators. Those high-level policies are sometimes error-prone and can fail to capture the dynamic nature of the applications. The results of these sub-optimal and laborious policies may lead to the under-utilization of physical clusters. This demands an optimal, dynamic resource allocation strategy which will not only meet the SLAs, but also will achieve the increased server utilization with minimal human intervention.

Many approaches have been proposed[3, 1] to address resource allocation in virtualized computing environments. These solutions broadly fall into two categories  either those came up with a model for a specific workload pattern (for example, database servers); or those strategies focused on specific subsystem (for example storage or memory). For the first category, modelers are required to know the details of the software architecture so that their solutions can take into account those specifics to increase the accuracy of their tools. Although, this is certainly feasible for a data center server running workloads of similar patterns, the approach does not fit for a server running applications with dissimilar resource usage characteristics. A generic strategy is required to cater to the need of diverse range of

applications. A recent study[2] came up with a model that is applicable for different workloads which may have multi-tier architectures. That study only focuses on controlling CPU shares. Still, its not clear how they will fit their model for all kinds of physical resources - CPU, memory, network and storage.

Our approach differs from the others in the sense that it tries to provide a holistic model which is applicable to all kinds of applications and which covers all the subsystems. In this paper, we suggest coarse-grained allocation policies for VMs without requiring much application details. We propose a model which is based on artificial neural network that treats each virtual machine as a black-box with the different allocation parameters as its input and the application performance metric as its output. The modeling employs two-steps. In the first step, individual applications are modeled by training each instance of Artificial Neural Network (ANN) with a set of observed application performance metric under different virtual machine configurations. The training will allow administrators to predict application output for a given set of input parameters. In normal conditions, the VMs are given resources as needed by the applications. In devising this model, we identified the minimal set of parameters that essentially incorporates all the factors that can influence the performance of an application running inside a VM. We also account the interferences (especially, I/O) caused by other VMs. Our set of parameters are general enough to be employed in other modeling techniques (e.g. regression analysis) also. We are going to discuss about the tuning parameters in details later in the paper. This entire part facilitates proper resource allocations to ensure the QoS as required by the strict SLAs. In the second phase, the model maps each application performance metric to the amount in dollars it contributes to the data center business operations. Our aim is to maximize the

profit across all the VMs. During contention for shared resources, we give extra resources to a VM with more shares by taking resources from a VM with fewer shares. Therefore, we will augment a VM that is more profitable from the point of view of the organization. In a nutshell, the model will automatically distribute the resources to a pool of VMs so that the total profit contributed by all the VMs is maximized.

We have implemented a prototype of our tool in Xen[4] Virtual Machine Monitor (VMM). For CPU, memory allocations we used the techniques already available in Xen VMM. We used other strategy which is based on I/O scheduler prioritization and contention from other virtual machines to control I/O for VMs. Initially, we experimented with postmark [?] which is an I/O intensive benchmark. Our ANN model with four parameters predicts the Transactions Per Second (TPS) of the postmark instance under different physical resource allocations to the VM running it within 5% accuracy. That gives us enough confidence to use this technique for modeling other applications as well.

The contributions of the paper are the followings: 1) we identified the parameters that are essential for accurate application modeling. 2) The design of a model that guarantees SLAs by distributing resources according to the proportional amount of shares. 3) The model also aims to maximize the overall profit for the data centers. Although, in this paper we are mainly concentrated in working on a single physical host, our approach can be easily extended into multiple physical clusters hosting several virtual machines.

# References

[1] A. Aboulnaga K. Salem P. Kokosielis A. A. Soror, U. F. Minhas and S. Ka-

math. Automatic virtual machine configuration for database workloads. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 953–966, 2008.

[2] P.Padala et. al. Adaptive control of virtualized resources in utility computing environments. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007*, pages 289–302, 2007.

[3] S. Ghanbari A. D. Popescu J. Chen G. Soundararajan, D. Lupei and C. Amza. Dynamic resource allocation for database servers running on virtual storage. In *Proceedings of FAST*, 2009.

[4] K. Fraser S. Hand T. Harris A. Ho R. Neugebauer P. Barham, B. Dragovic. Xen and the art of virtualization. In *Proceedings of the nineteenth ACM symposium on Operating systems principles*, pages 164–177, 2003.