# Multidimensional Resource Consumption Analysis of Co-Located VMs using PCA

Waltenegus Dargie

Chair for Computer Networks, Faculty of Computer Science, Technical University of Dresden, 01062 Dresden, Germany Email:waltenegus.dargie@tu-dresden.de

Abstract—One of the strategies employed to deal with resource inefficiency in data centres is dynamic virtual machine/container consolidation. The idea behind is, by populating physical servers with an optimal number of virtual machines, all the server's resources (CPU, memory, network bandwidth, etc.) can be utilised effectively. This approach requires (1) the free migration of virtual machine at runtime and (2) the identification of virtual machines which exhibit complementary features. Most existing or proposed approaches are based on elaborate and complex multi-variate optimisation and do not easily lend themselves to fast and intuitive solutions. In this paper, we investigate the scope and usefulness of dimensionality reduction techniques, ideas borrowed from unsupervised machine learning, to analyse the existence of contentious and complementary features in the resource consumption characteristics of co-located virtual machines. Initial results suggest that fast and tractable scheduling can be achieved using these techniques.

Index Terms—Energy-efficient computing, principal component analysis, PCA, scheduling, virtual machine consolidation

### I. INTRODUCTION

The introduction of virtualization and cloud computing has enabled the relative efficient use of computing resources in data centres [1], [2]. At present, a large amount of virtual machines/containers<sup>1</sup> belonging to different owners can be hosted by a single physical machine. Moreover, the VMs can be seamlessly migrated at runtime from idle or underutilised servers to other host servers, so that the formers can be switched-off to save power [3]. However, aggregating VMs onto a few number of servers may lead to overloading the target servers resulting in a disproportional amount of power consumption and the violation of service level agreements on account of performance degradation.

Managing computing resources at runtime is often an NPhard optimisation problem [4], as the task should take into consideration many factors, including, the resource consumption characteristics of the VMs being managed, the available network bandwidth for migration as well as the data centre topology. The resource consumption characteristic itself is a multi-dimensional aspect [5]. Consequently, fully utilising all the resources (CPU, cache, memory, network bandwidth, etc.) of a physical server is a difficult objective to achieve even for elaborate and complex approaches.



Fig. 1: The normalised workload of YouTube for April 10 and 11 2018. The normalised workload is computed as the ratio of the incoming requests of YouTube Germany to the requests placed at the same time to all YouTube servers worldwide<sup>2</sup>.

Ideally, a consolidation strategy should co-locate VMs having complementary resource utilisation characteristics, complementary both in terms of the resources they utilise at any given time (for example, CPU-intensive and IO-intensive) and in time (having non-overlapping executing time). Different experimental studies suggest that even though the short-term workloads of typical Internet applications should be regarded as random processes, their long-term resource consumption characteristics, nevertheless, have predictable patters [6], [7]. This can be illustrated by examining the workload pattern of YouTube Germany (ref. to Fig. 1). As can be seen, the shortterm workload fluctuated considerably between April 10 and 11, 2018, but when we regard the day-to-day workload, the pattern remained more or less unchanged. Our examination of 16 different Google applications enables us to conclude that this characteristic ubiquitous among typical Internet applications.

A significant portion of a consolidation assignment can be carried out offline, taking into account the long-term resource consumption statistics of hosted VMs. In this paper, we analyse the resource consumption statistics of 36 real-world VMs which are currently running on our Chair's main server at the Faculty of Computer Science. Our objective is to examine the existence of multidimensional contentious and complementary aspects and how this knowledge can be exploited to efficiently consolidate the VMs.

The remaining part of the paper is organized as follows: In section II we briefly summarise how we acquired the

<sup>2</sup>https://transparencyreport.google.com/traffic/overview

<sup>&</sup>lt;sup>1</sup>Henceforth we collectively call them VMs for convenience.

measurements sets for our analysis. In section III, we give an overview of dimensionality reduction techniques with the focus on Principal Component Analysis (PCA) and analyse the existence of contention and complementary features in resource utilisation when multiple virtual machines are executed on one and the same physical server. Finally, in section IV, we provide concluding remarks and outline future work.

# II. BACKGROUND

For the analysis we present in this paper, we rely on statistics obtained from the main server of the Chair of Computer Networks at the TU Dresden. The server currently hosts 36 active virtual machines belonging to different owners. We do not have knowledge of the specific purpose of the VMs but were able to synchronously sample their resource utilisation. Our analysis is based on the three-hour statistics we gathered this way. Even though our measurement sets are extensive, we limit our analysis to the utilisation of the three most important resources, namely, the CPU, the memory (MEM), and the IO (NET), in order to make our analysis visually tractable. The server is a PRIMERGY RX300 S6, Fujitsu rack server, optimised for virtualization purposes. It has two physical sockets housing two 2.4 GHz Intel Xeon E5-620 quad-core processors. When in a Hyper-Threading mode, the server provides 16 logical cores altogether. In addition, it has 48 GB DDR3 RAM memory and runs VMware vSphere Hypervisor  $5.5.0^3$ .

Fig. 2 displays the standardised CPU and MEM utilisation of two of the virtual machines. As can be seen, the two VMs are markedly different in their MEM utilisation whereas their CPU utilisation pattern is more or less the same. The figure highlights the difficulty of achieving complementary consolidation. For example, if we co-locate these VMs, we can achieve complementarity in terms of memory utilisation, but not in terms of CPU utilisation. This clearly highlights the limitation of single dimensional VM consolidation strategies [8].

# III. PRINCIPAL COMPONENT ANALYSIS

Considering the large amount of data one has to process during multidimensional optimisation, most statistical approaches are intractable even for offline consolidation. This is particularly true when there is a strong statistical dependence between the resources consumed by a single virtual machine as well as between multiple virtual machines<sup>4</sup>. For example, we have observed a strong correlation between the memory utilisation and the network bandwidth utilisation in some of our virtual machines, as these two resources coordinated when dealing with Internet traffic.

<sup>3</sup>https://www.vmware.com/de.html.

<sup>4</sup>Suppose we decide to co-locate two virtual machines based on their CPU utilisation statistics. The overall CPU utilisation is determined by convolving the pdfs of the CPU utilisation of the two VMs. If the VMs are statistically independent, the convolution operation yields a relatively tractable solution, because the joint pdf can be expressed as the multiplication of the individual pdfs. If, however, the VMs are not statistically independent, the convolution operation is in general difficult to solve.



Fig. 2: The standardised time series of the CPU and MEM utlisation of two different virtual machines.

TABLE I: A snapshot of the standardised resource consumption of VM1.

Sample	CPU	MEM	NET
1	0.7597341	-0.562566	1.4366341
2	1.6617245	3.237548	-0.7995181
3	:	:	:

Dimensionality reduction techniques, such as Principal component analysis (PCA) [9], [10] and tensor decomposition [11], are ideal complements to pdf based multidimensional optimisations when statistical independence between the VMs as well as their resource consumption cannot be guaranteed. These approaches, by identifying statistical dependencies, transform a high-dimensional dataset into a tractable, low-dimensional form without losing much information. We employed PCA to investigate statistical dependencies between VMs. One of the advantages of using a PCA is the flexibility of organising and interpreting the measurement sets.

## A. Dependency within a Single VM

When analysing the resource consumption characteristics of a single VM, the measurement sets can be represented by a matrix. Depending on which parameters are regarded as the rows and which as the columns of the matrix, different objectives can be achieved.



TABLE II: The resource consumption of VM1 expressed by two principal components.

PC2

PC1

Samples

Fig. 3: The resource consumption of four of the most active VMs expressed by two principal components.

- When the resources are arranged as the columns of the matrix and the samples as the rows (ref. to Tab. I), PCA reduces the dimensionality of the matrix by exploiting dependency between the resources consumed but leaves the sample size of the matrix (the rows) intact (ref. to Tab. II). In other words, the temporal characteristic of the resource consumption is preserved but explicit understanding, whether or not a VM is compute-intensive, data-intensive, or communication intensive, cannot be made.
- When, on the other hand, the samples make up the columns of the matrix and the resources make up the rows, PCA can be employed to exploit periodic features in order to significantly reduce the dimension of the dataset. However, the number of resources considered in the investigation remains intact. The advantage of this approach is that now the VMs can be classified as compute-intensive, data-intensive, or communication intensive VMs.

Fig. 3 displays the bi-plots of the resource utilisation (CPU, MEM, NET) of four virtual machines explained by two principal components with 99% accuracy. The measurement

TABLE III: A snapshot of the standardised CPU utilisation of all the VMs.

	1	2	3	4	5	6		
V1	0.76	1.67	-0.97	-0.63	-1.0	1.11		
V2	-0.67	0.066	0.98	-0.69	-0.70	1.99		
V3	:		:	÷		:	:	:

sets were initially organised and preprocessed (centred and scaled) as in Tab. I. The significance of the four axes in each plot is explained as follows: the horizontal (bottom) and the vertical (left) axes depict the first and the second principal components, respectively. The horizontal (top) and the vertical (right) axes depict the scores of the samples and the loadings of the original variables<sup>5</sup>, respectively. The red arrows indicate in which direction and with what magnitude each resource utilisation varies with respect to the principal components.

If we compare the virtual machines in terms of the relative direction of the resources, we can identify complementary and contentious features. For example, the first virtual machine (the plot in the second quadrant) and the second (the plot in the first quadrant) have almost orthogonal relationships. The memory utilisation of the first virtual machine contributes almost nothing to the first principal component whereas its contribution to the second principal component is relatively large and positive. Whereas for the second virtual machine, the contribution of the memory utilisation to the first principal component is large and negative. A similar comparison can be made in terms of the network bandwidth utilisation. This implies that scheduling these two virtual machines to execute on the same physical server will result in an efficient resource utilisation. The same can be said of the second and the third (third quadrant) virtual machines.

By comparison, consolidating the first and the fourth (fourth quadrant) virtual machines will result in a significant resource contention. As can be seen, the magnitude and direction of the contributions of the three resources to the principal components are similar suggesting that the two virtual machines have similar resource consumption characteristics.

The 2-dimensional VM-by-VM analysis is comprehensible as long as the number of hosted VMs is small. When the number becomes sizable, however, making an objective comparison becomes difficult.

# B. Dependency Between All VMS

A more comprehensible clustering can be achieved by analysing the utilisation of each resource by all the VMs. For instance, Fig. 4 displays the CPU utilisation of all the VMs explained by three principal components. For this analysis, we formed a  $36 \times 10800$  matrix by putting together the CPU

<sup>&</sup>lt;sup>5</sup>Formally, the PCA is based on a decomposition of the data matrix **A** into two matrices **V** and **U**:  $\mathbf{A} = \mathbf{U}\mathbf{V}^{\mathsf{T}}$ . **V** and **U** are orthogonal matrices. The former is called the loading matrix whereas the latter is called the scores matrix. The loadings can be understood as the weights for each original variable (for our case, CPU, MEM, and NET) when calculating the principal components. The matrix **U** contains the original data in a rotated coordinate system.

TABLE IV: The CPU utilisation of all VMs represented by three principal components.

	PC1	PC2	PC3
VM1	-4.733735	1.21015794	-0.3178977
VM2	7.116493	-3.56787504	-6.5977242
VM3	-3.063336	-0.07835434	1.0430113
:	:	:	:



Fig. 4: The CPU utilisation of all the VMs expressed by three principal components.

utilisation of all the VMS, the rows of the matrix constituting the VMs and the columns of the matrix constituting the sample instances.

Fig. 4 displays the relative CPU usage of all the VMs using three principal components. Each dot in the plot represents a VM. Most of the VMS are clustered towards the origin, indicating that their CPU utilisation does not exhibit a remarkable variation throughout. These VMs can be scheduled together provided that their mean CPU utilisation is small. Alternatively, it can be said that the CPU utilisation of these VMs can be regarded, by and large, as deterministic and easily lends itself to deterministic scheduling. On the other hand, the dots displayed towards the top and the bottom of the box as well as away from the origin indicate the VMs whose CPU utilisation shows high variability, implying that these VMs, when scheduled together, may result in great fluctuations of resource consumption leading to either considerable underutilisation or overloading of the server. By contrast, scheduling VMs on the opposite sides of the box results in appreciable complementarity.

#### **IV. DISCUSSION**

In this paper we established the groundwork for a scalable, multi-dimensional virtual machine consolidation in large-scale data centres. We investigated the existence of correlation and anti-correlation between co-located virtual machines using principal component analysis (PCA), a dimensionality reduction technique which has a wide range of applications in big data analytics. In addition to its capacity to significantly reduce the dimension of an observation (measurement sets), PCA also enables to cluster virtual machines, so that contentious and complementary features can be identified and variability analysis can be carried out. Our analysis was focused on the CPU, MEM, and network bandwidth utilisation of 36 realworld virtual machines.

Initial results suggest that using only 10 principal components, the existence of complementarity and contention amongst co-located virtual machines can easily be determined with 70 % accuracy for CPU-intensive virtual machines. This amounts to reducing a  $36 \times 10800$  measurement set to  $36 \times 10$ measurement sets. Likewise, the analysis on the memory utilisation required just 100 principal components to achieve more than 90 % accuracy to identify contention, thereby reducing the complexity of the analysis task by nearly 99 %.

In future we plan to extend our work to tensor decomposition, so that the resources consumption of all the virtual machines can be represented as a three dimensional array and can be analysed in a single step. This way a more comprehensive insight can be gained pertaining to the statistical dependence between the hosted virtual machines.

#### ACKNOWLEDGMENT

This work has been partially funded by the German Research Foundation (DFG) under project agreement: SFB 912: Highly Adaptive Energy-Efficient Computing (HAEC).

#### REFERENCES

- W. Dargie, "A stochastic model for estimating the power consumption of a processor," *IEEE Transactions on Computers*, vol. 64, no. 5, pp. 1311– 1322, 2015.
- [2] C. Mobius, W. Dargie, and A. Schill, "Power consumption estimation models for processors, virtual machines, and servers," *IEEE Transactions* on Parallel and Distributed Systems, vol. 25, no. 6, pp. 1600–1614, 2014.
- [3] A. Strunk and W. Dargie, "Does live migration of virtual machines cost energy?," in Advanced Information Networking and Applications (AINA), 2013 IEEE 27th International Conference on, pp. 514–521, IEEE, 2013.
- [4] M. H. Ferdaus, M. Murshed, R. N. Calheiros, and R. Buyya, "Virtual machine consolidation in cloud data centers using aco metaheuristic," in *European Conference on Parallel Processing*, pp. 306–317, Springer, 2014.
- [5] A. Ashraf and I. Porres, "Multi-objective dynamic virtual machine consolidation in the cloud using ant colony system," *International Journal of Parallel, Emergent and Distributed Systems*, vol. 33, no. 1, pp. 103–120, 2018.
- [6] R. N. Calheiros, E. Masoumi, R. Ranjan, and R. Buyya, "Workload prediction using arima model and its impact on cloud applications' qos," *IEEE Transactions on Cloud Computing*, vol. 3, no. 4, pp. 449–458, 2015.
- [7] M. Haehnel, J. Martinovic, G. Scheithauer, A. Fischer, A. Schill, and W. Dargie, "Extending the cutting stock problem for consolidating services with stochastic workloads," *IEEE Transactions on Parallel and Distributed Systems*, pp. 1–1, 2018.
- [8] S. Wang, A. Zhou, C.-H. Hsu, X. Xiao, and F. Yang, "Provision of dataintensive services through energy-and qos-aware virtual machine placement in national cloud data centers," *IEEE Transactions on Emerging Topics in Computing*, vol. 4, no. 2, pp. 290–300, 2016.
- [9] R. Bro and A. K. Smilde, "Principal component analysis," Analytical Methods, vol. 6, no. 9, pp. 2812–2831, 2014.
- [10] J. Shlens, "A tutorial on principal component analysis," arXiv preprint arXiv:1404.1100, 2014.
- [11] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," SIAM review, vol. 51, no. 3, pp. 455–500, 2009.