ELSEVIER

CrossMark

# An algorithm for network and data-aware placement of multi-tier applications in cloud data centers

Md Hasanul Ferdaus[a,b,*], Manzur Murshed[c], Rodrigo N. Calheiros[d], Rajkumar Buyya[b]

[a] Faculty of Information Technology, 25 Exhibition Walk, Clayton campus, Monash University, VIC 3800, Australia
[b] Cloud Computing and Distributed Systems (CLOUDS) Lab, School of Computing and Information Systems, The University of Melbourne, Australia
[c] Faculty of Science and Technology, Federation University Australia, Northways Road, Churchill, VIC 3842, Australia
[d] School of Computing, Engineering and Mathematics, Western Sydney University, Australia

## ARTICLE INFO

## ABSTRACT

Today's Cloud applications are dominated by composite applications comprising multiple computing and data components with strong communication correlations among them. Although Cloud providers are deploying large number of computing and storage devices to address the ever increasing demand for computing and storage resources, network resource demands are emerging as one of the key areas of performance bottleneck. This paper addresses network-aware placement of virtual components (computing and data) of multi-tier applications in data centers and formally defines the placement as an optimization problem. The simultaneous placement of Virtual Machines and data blocks aims at reducing the network overhead of the data center network infrastructure. A greedy heuristic is proposed for the on-demand application components placement that localizes network traffic in the data center interconnect. Such optimization helps reducing communication overhead in upper layer network switches that will eventually reduce the overall traffic volume across the data center. This, in turn, will help reducing packet transmission delay, increasing network performance, and minimizing the energy consumption of network components. Experimental results demonstrate performance superiority of the proposed algorithm over other approaches where it outperforms the state-of-the-art network-aware application placement algorithm across all performance metrics by reducing the average network cost up to 67% and network usage at core switches up to 84%, as well as increasing the average number of application deployments up to 18%.

## 1. Introduction

With the pragmatic realization of computing as a utility, Cloud Computing has recently emerged as a highly successful alternative information technology paradigm through the unique features of on-demand resource provisioning, pay-as-you-go business model, virtually unlimited amount of computing resources, and high reliability (Buyya et al., 2009). In order to meet the rapidly increasing demand for computing, communication, and storage resources, Cloud providers are deploying large-scale data centers comprising thousands of servers across the planet. These data centers are experiencing sharp rise in network traffic and a major portion of this traffic is constituted of the data communication within the data center. Recent report from Cisco Systems Inc. (Cisco, 2015) demonstrates that the Cloud data centers will dominate the global data center traffic flow for the foreseeable future and

its importance is highlighted by one of the top-line projections from this forecast that, by 2019, more than four-fifths of the total data center traffic will be Cloud traffic (Fig. 1). One important trait pointed out by the report is that a majority of the global data center traffic is generated due to the data communication within the data centers: in 2014, it was 75.4% and it will be around 73.1% in 2019.

This huge amount of intra-data center traffic is primarily generated by the application components that are correlated to each other, for example, the computing components of a composite application (e.g., MapReduce) writing data to the storage array after it has processed the data. This large growth of data center traffic may pose serious scalability problems for wide adoption of Cloud Computing. Moreover, by the way of continuously rising popularity of social networking sites, e-commerce, and Internet-based gaming applications, large amount of data processing has become an integral part of Cloud applications. Furthermore,
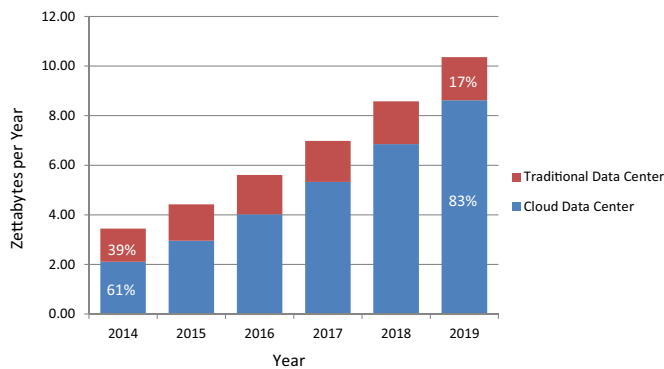
---

**Fig. 1.** Worldwide data center traffic growth (data source: Cisco).

scientific processing, multimedia rendering, workflow, and other massive parallel processing and business applications are being migrated to the Clouds due to the unique advantages of high scalability, reliability, and pay-per-use business model. Over and above, recent trend in Big Data computing using Cloud resources (Assuncao et al., 2015) is emerging as a rapidly growing factor contributing to the rise of network traffic in Cloud data centers.

One of the key technological elements that have paved the way for the extreme success of Cloud Computing is virtualization. Modern data centers leverage various virtualization technologies (e.g., machine, network, and storage virtualization) to provide users an abstraction layer that delivers a uniform and seamless computing platform by hiding the underlying hardware heterogeneity, geographic boundaries, and internal management complexities (Zhang et al., 2010). By the use of virtualization, physical server resources are abstracted and shared through partial or full machine simulation by time-sharing, and hardware and software partitioning into multiple execution environments, known as *Virtual Machines* (VMs), each of which runs as a complete and isolated system. It allows dynamic sharing and reconfiguration of physical resources in Cloud infrastructures that make it possible to run multiple applications in separate VMs having different performance metrics. It also facilitates Cloud providers to improve utilization of physical servers through VM multiplexing (Meng et al., 2010a) and multi-tenancy, i.e., simultaneous sharing of physical resources of the same server by multiple Cloud customers. Furthermore, it enables on-demand resource pooling through which computing (e.g., CPU and memory), network, and storage resources are provisioned to customers only when needed (Kusic et al., 2009). By utilizing these flexible features of virtualization for provisioning physical resources, the scalability of data center network can be improved through minimization of network load imposed due to the deployment of customer applications.

On the other side, modern Cloud applications are dominated by multi-component applications such as multi-tier applications, massive parallel processing applications, scientific and business workflows, content delivery networks, and so on. These applications usually have multiple computing and associated data components. The computing components are usually delivered to customers in the form of VMs, such as Amazon EC2 Instances,[1] whereas the data components are delivered as data blocks, such as Amazon EBS.[2] These computing components of such applications have specific service roles and are arranged in layers in the overall structural design of the application. For example, large enterprise applications are often modeled as 3-tier applications: the presentation tier (e.g., web server), the logic tier (e.g., application server), and the data tier (e.g., relational database) (Urgaonkar et al., 2005). The computing components (VMs) of such applications have specific communication requirements among them-

selves, as well as with the data blocks that are associated to those VMs (Fig. 2). As a consequence, overall performance of such applications highly depends on the communication delays among the computing and data components. From the Cloud providers' perspective, optimization of network utilization of data center resources is tantamount to profit maximization. Moreover, efficient bandwidth allocation and reduction of data packet hopping through network devices (e.g., switches or routers) trim down the overall energy consumption of network infrastructure. On the other hand, Cloud consumers' concern is to receive guaranteed Quality of Service (QoS) of the delivered virtual resources, which can be assured through appropriate provisioning of requested resources.

Given the issues of sharp rise in network traffic in data centers, this paper addresses the scalability concern of data center network through a traffic-aware placement strategy of multi-component, composite application (in particular, VMs and data blocks) in virtualized data center that aims at optimizing the network traffic load incurred due to placement decision. Such placement decisions can be made during the application deployment phase in the data center. VM placement decisions focusing on other goals rather than network efficiency, such as energy consumption reduction (Feller et al., 2011; Beloglazov and Buyya, 2012) and server resource utilization (Gao et al., 2013; Ferdaus et al., 2014), often result in placements where VMs with high mutual traffic are placed in host servers with high mutual network cost. For example, one of our previous works (Ferdaus et al., 2014) on the placement of a cluster of VMs strives to consolidate the VMs into a minimal number of servers in order to reduce server resource wastage. By this process, unused servers can be kept into lower power states (e.g., suspended) so as to improve power efficiency of the data center. Since this approach does not consider inter-VM network communication patterns, such placement decisions can eventually result in locating VMs with high mutual network traffic in long distant servers, such as servers locating across the network edges. Several other VM placement works focusing on non-network objectives can be found in (Wu and Ishikawa, 2015; Farahnakian et al., 2015; Nguyen et al., 2014; Corradi et al., 2014; Alboaneen et al., 2014). With a network-focused analysis, it can be concluded that research works such as the above ones considered single-tier applications and VM clusters without consideration of mutual network communication within the application components or VMs. On the contrary, this paper focuses on placing mutually communicating components of applications (such as VMs and data blocks) in data center components (such as physical servers and storage devices) with lesser network cost so that network overhead imposed due to the application placement is minimized. With this placement goal, the best placement for two communicating VMs would be in the same server where they can communicate through memory copy, rather than using the physical network links. This paper effectively addresses network-focused placement problem of multi-tiered applications with components having mutual network communication rather than single-tiered ones. The significance of the network-focused placement of multi-tiered applications is evident from the experimental results presented later in Section 5, where it is observed that an efficient non-network greedy placement algorithm, namely First Fit Decreasing (FFD), incurs higher network costs compared to the proposed network-aware placement heuristic.

Moreover, advanced hardware devices with combined capabilities are opening new opportunities for efficient resource allocation focusing on application needs. For example, Dell PowerEdge C8000 moduler servers are equipped with CPU, GPU, and storage components that can work as multi-function devices. Combined placement of application components with high mutual traffic (e.g., VMs and their associated data components) in such multi-function servers will effectively reduce the data transfer delay since the data accessed by the VMs reside in the same devices. Similar trends are found in high-end network switches (e.g., Cisco MDS 9200 Multiservice Switches) that come with additional built-in processing and storage capabilities. Reflecting on these tech-

---

[1] Amazon EC2 - Virtual Server Hosting, 2016. https://aws.amazon.com/ec2/.
[2] Amazon Elastic Block Store (EBS), 2016. https://aws.amazon.com/ebs/.