

# IoNCloud: approaching high utilization and network predictability with shared bandwidth guarantees

Miguel C. Neves, Daniel S. Marcon, Rodrigo R. Oliveira, Tiago de Almeida,  
Leonardo R. Bays, Luciano P. Gaspary, Marinho P. Barcellos

Institute of Informatics – Federal University of Rio Grande do Sul

Email: {mcneves, daniel.stefani, ruas.oliveira, talmeida, lrbays, paschoal, marinho}@inf.ufrgs.br

**Abstract**—The intra-cloud network is shared in a best-effort manner, which causes tenant applications to have no actual bandwidth guarantees. In this paper, we introduce a resource allocation strategy that aims at minimizing resource underutilization with low management overhead. To demonstrate the benefits of the strategy, we develop IoNCloud, a system that implements the proposed allocation scheme. IoNCloud employs the concept of ion attraction/repulsion to group applications in virtual networks according to their temporal bandwidth demands. In doing so, we explore the trade-off between high resource utilization (desired by providers) and strict network guarantees (desired by tenants).

## I. INTRODUCTION

Cloud computing allows tenants to execute several kinds of applications, both inward computation with bandwidth-intensive network usage and user-facing ones with strict response times. Providers, however, offer no actual bandwidth guarantees for applications [1], [2]. The intra-cloud network is typically oversubscribed (more bandwidth available in leaf nodes than in the core) and is shared in a best-effort manner, relying on TCP to achieve high utilization. TCP, nonetheless, does not provide robust isolation among flows<sup>1</sup> in the network; in fact, long-lived flows with a large number of packets are privileged over small ones (which is typically called *performance interference*) [2]. Moreover, recent studies [3] show that bandwidth available for virtual machines (VMs) in the intra-cloud network can vary by a factor of five or more, resulting in poor and unpredictable overall application performance.

## II. IONCLOUD

We propose a resource allocation strategy for reserving and isolating network resources in cloud datacenters. It aims at minimizing resource underutilization while providing an efficient way to predictably share bandwidth among applications. To show the benefits of the strategy, we develop IoNCloud (Isolation of Networks in the Cloud), a system that implements the proposed allocation scheme. IoNCloud groups tenant applications into virtual networks (VNs) with bandwidth guarantees, according to their temporal network usage and need of isolation. In doing so, we seek to explore the trade-off between high resource utilization (a common goal for providers to reduce operational costs and achieve economies of scale) and strict network guarantees (desired by tenants).

In this strategy, all applications that belong to the same group share the same set of (virtual) network resources (i.e., they have shared bandwidth guarantees). In contrast, virtual networks are completely isolated from one another,

<sup>1</sup>Flows are characterized by sequences of packets from a source to a destination host.

which means that each group has a guaranteed amount of network resources. An abstract view of IoNCloud is shown in Figure 1, which depicts application requests being received and allocated in two steps. The first step is responsible for application demand analysis and grouping in VNs, while the second embeds virtual networks onto the physical substrate.

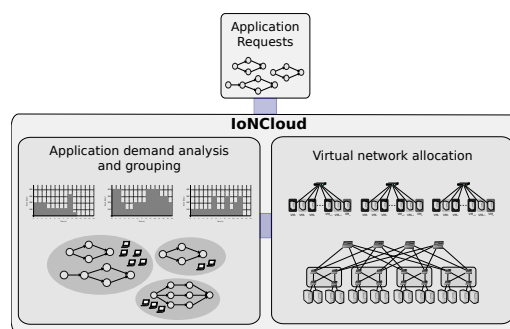


Fig. 1: IoNCloud system overview.

## III. CONCLUSION AND FUTURE WORK

In general, the proposed approach introduces a network-performance-aware resource allocation strategy for Infrastructure as a Service (IaaS) cloud platforms. It improves network predictability by grouping tenant applications into virtual networks according to their temporal bandwidth demands. Furthermore, we develop IoNCloud, a system that implements the proposed strategy for large-scale cloud datacenters. IoNCloud groups applications in VNs, maps them on the physical substrate and provisions network resources according to the aggregate temporal bandwidth demands of the applications in each group. In future work, we intend to perform an extensive evaluation of the system as well as to add VM migration to minimize network traffic.

## ACKNOWLEDGEMENTS

This work has been supported by FP7/CNPq (Project SecFuNet, FP7-ICT-2011-EU-Brazil), RNP-CTIC (Project AltoStratus), as well as PRONEM/FAPERGS/CNPq (Project NPRV).

## REFERENCES

- [1] D. Xie, N. Ding, Y. C. Hu, and R. Kompella, “The only constant is change: Incorporating time-varying network reservations in data centers,” in *ACM SIGCOMM*, 2012.
- [2] D. Abts and B. Felderman, “A guided tour of data-center networking,” *Comm. ACM*, vol. 55, no. 6, Jun. 2012.
- [3] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron, “Towards predictable datacenter networks,” in *ACM SIGCOMM*, 2011.