

Impact of Network Function Virtualization: A Study based on Real-Life Mobile Network Data

Ahmad Bilal*, Andras Vajda* and Taleb Tarik†

* Ericsson, Finland

† Aalto University, Finland

bilal.x.ahmad@ericsson.com, andras.vajda@ericsson.com, talebtarik@ieee.org

Abstract—Mobile Operators are looking for new ways to cope with ever-increasing data traffic while improving the operational and capital efficiency of their networks. Cloud computing and network function virtualization (NFV) have emerged as key enablers to optimize resource utilization and at the same time reduce network operational expenditure (OPEX). In virtualized networks, network functions are delivered as software running on generic hardware allowing service providers to dynamically allocate resources based on traffic and service demands. In this paper, we analyze resource utilization using real-life data of two different mobile networks and evaluate the impact virtualization would have on these networks. Some conclusions are drawn based on the analysis.

Index Terms—NFV, Cloud Computing, Real-Life Mobile Network Data, Resource utilization.

I. INTRODUCTION

Traditionally, network nodes are delivered pre-configured in a highly optimized manner with specialized hardware specific to node functionality. Deployment of new network services typically requires separate hardware with significant cost and lead-time of integration and operation. Network Function Virtualization (NFV) [1] promises to address these challenges through the decoupling of software from hardware with the introduction of a virtualization layer. NFV can be applied to most functions in the network whether these are control-plane or data-plane functions. The goal of NFV is to run network functions as software in e.g. virtual machines (VMs) on top of virtualization platforms deployed on generic hardware. It is understood that running network functions on general purpose hardware, rather than on dedicated hardware, can impact performance. The virtualization layer introduces latency and an extra overhead that consumes extra capacity. Significant overhead may be required to implement software-based switches (commonly called virtual switches or vSwitch) that route packets to and from appropriate VMs. This CPU overhead can reduce maximum throughput and increase latency on an I/O device.

In this paper, real data from two different packet core networks with dedicated network nodes were analyzed. The two selected packet core networks represent mobile network deployments in developed and developing countries, respectively. We created a load profile for each node that allowed us to evaluate how the overall load evolves over a period of time. Based on the load profile, the amount of general purpose hardware needed was calculated, assuming that the

same type of load is executed in a virtualized environment. We also analyzed different scenarios with various virtualization overheads for control and data-plane processing. To the best knowledge of the authors, this is the first measurement study of live network resources utilization which will give future directions to implementation of NFV in mobile networks.

The remainder of this paper is organized as follows. Section II presents background literature related to NFV technologies. Section III presents the networks and key nodes under study. Section IV reports our analysis based on the measured data. The simulation results with and without virtualization overheads for both networks are also reported therein. Finally, the paper concludes in Section V.

II. NFV TECHNOLOGIES: RELATED WORK

NFV allows a service provider to deliver network functions as pure software running in a virtualized environment with reduced cost and high deployment efficiency [6][7]. This shift of hardware to software running in standard virtual machines or containers (e.g., Docker or Googles Kubernetes) is expected to reduce CAPEX and OPEX. Reduced cost, increased service deployment velocity, services introduced based on geography and costumers' needs, ability to efficiently cope with emerging resource-intensive applications [8], reduced energy consumption and several other benefits are expected to be achieved through NFV [1]. In short, NFV, along with Software Defined Networking, will enable the launch of anything as a service in a more cost efficiency way while ensuring short time of service to market [9]. However, this cannot be achieved without addressing challenges spanning from system design to ensuring service resiliency [10].

NFV is also foreseen as an important technology to enable the on-demand creation of cloud-based virtual mobile networks. Here, an important challenge pertains to the placement of Virtualized Network Functions (VNFs), within the same or across distributed datacenters, considering the performance constraints and functional relationship among VNFs that form a single virtual network infrastructure. To this problem, different solutions have been devised not only ensuring communications efficiency for mobile users placing VNFs at strategic positions but also ensuring cost efficiency for the operators; minimizing the cost associated with the instantiated VMs. The work in [11] analyzes the impact of deployment strategy on the overall performance of virtual network infrastructures

deployed on the cloud. In [12], a fine granular resource-aware VNF management is proposed for the initial deployment and runtime management of virtual network infrastructures.

Cloud computing and industry-standard high-volume servers are key enablers to achieve the goals of NFV. However, general purpose processors are not designed to process modern, high speed protocols [2]. Historically, control plane applications have been implemented on general purpose hardware and the performance degradation is small enough to be neglected. On the other hand, data plane applications are mostly executed on specialized hardware to perform specific functions and to meet specific requirements in terms of Quality of Experience (QoE) or Quality of Service (QoS). Therefore, realizing data plane application on general purpose hardware to meet desired performance and cost requirements is an industry challenge that needs to be tackled. In NFV, performance degradation is due to the introduction of a virtualization layer between hardware and Operating System (OS). This abstraction layer, i.e., hypervisor, increases the overhead for accessing resources, most importantly I/O resources, which are the most difficult to partition and manage [3], due to mismatch between CPU and I/O speeds; random nature of packet arrival; higher number of virtual machines that need to be served; and overhead due to switching done purely in software (i.e., the vSwitch component of hypervisor).

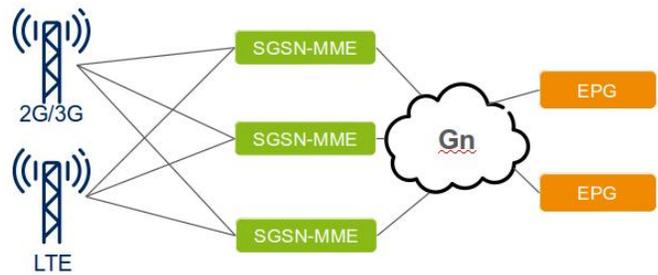
There are software-based, hardware-assisted and full hardware-based solutions for sharing I/O devices. With the software-based solutions, the hypervisor translates the I/O requests it receives and allocates them serially to physical resources and packets are routed through the virtualization layer, which adds performance overhead. Virtual machine performance is greatly degraded if the processor is interrupted frequently for I/O [4].

Some network devices use hardware-assisted approach to accelerate software-based I/O sharing but this solution still involves a hypervisor in between, resulting in performance degradation. In order to overcome this performance overhead the solution needs to avoid the involvement of hypervisors in I/O operations. The PCI-SIG SR-IOV standard [5] offers an approach for I/O sharing that overcomes the performance limitations as discussed above. SR-IOV enables I/O devices to support multiple virtual functions, lightweight mechanisms for transferring data to and from network adapter, which can be assigned directly and can communicate to VM bypassing hypervisor [4]. Although this solution is highly efficient in terms of performance, it has important drawbacks [3], in terms of scalability, outbound traffic and live migration limitations.

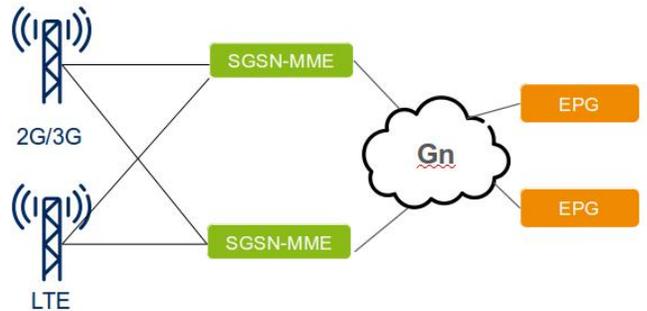
Besides all these solutions, there is no study on real-life networks that can quantify the impact of virtualization on overall system capacity and performance if virtualization and data plane overhead are taken into account.

III. ANALYZED NETWORKS

This section provides an overview of the studied mobile networks. The main nodes within the packet core network under study are:



(a) West European country network (NW 1)



(b) African country network (NW 2)

Fig. 1: The two studied mobile networks

- 1) SGSN-MME (Serving GPRS Support Node - Mobility Management Entity): The SGSN-MME provides SGSN and MME functionality. It is possible to use SGSN-MME with SGSN functionality, MME functionality or both.
- 2) EPG (Evolved Packet Gateway): The EPG provides GGSN (Gateway GPRS Support Node), S-GW (Serving Gateway), and P-GW (Packet Data Network Gateway) functionality. It is possible to use EPG with GGSN functionality, S-GW functionality, P-GW functionality, S-GW and P-GW functionality, or all simultaneously.

A. West European Country Network (NW 1)

There are three packet core sites in three different locations consisting of three SGSN-MMEs and two EPGs as shown in Fig. 1(a). The three SGSN-MMEs are integrated in a triple access pool network that provides redundancy to packet-switched network. 2G and 3G traffic from BSCs (Base Station Controller) and RNCs (Radio Network Controller) are distributed in different ratios among three geographical areas. 4G traffic is disturbed equally among the SGSN-MMEs. The two EPGs are also configured for triple access (2G/3G/LTE) at different locations. Each EPG is divided into two logical nodes: S-GW is only used for 4G traffic, P-GW used for 2G/3G and 4G connections. For NW 1, we used high granularity data over four months for SGSN-MME resources and three months for EPG nodes.

B. African Country Network (NW 2)

There are two packet core sites in two different locations consisting of two SGSN-MMEs and two EPGs as shown in Fig. 1(b). This network also provides packet services to 2G/3G and 4G users. However, resources utilization is small compared to NW1 due to low number of PS users. Furthermore, there is difference in resources dedicated to each network; for instance the processor resources in NW 2 are one third of those in NW 1. For this network, SGSN-MMEs data is analyzed for seven months while EPG data is analyzed for five months with one month excluded in between. The second EPG is not in use. For the sake of resiliency, it is used as back up if running EPG goes down.

The motivation behind studying these two networks is to compare cost-effectiveness between virtualizing a network with higher density of packet services and a network with low density of packet services. For both networks, the SGSN-MMEs statistics are captured with a five minutes granularity while EPG statistics are captured with a fifteen minutes granularity.

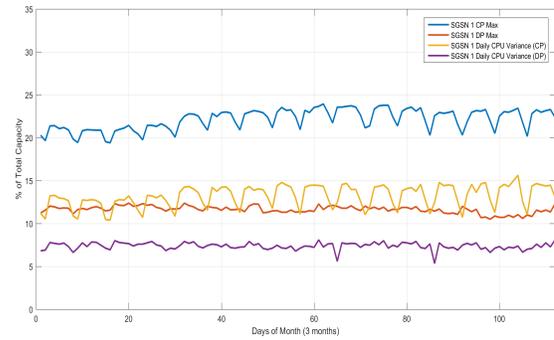
IV. ANALYSIS

This section is divided into two parts. The first part introduces the actual measured data while the second part shows simulation results when the actual measured data was to be executed in a virtual environment.

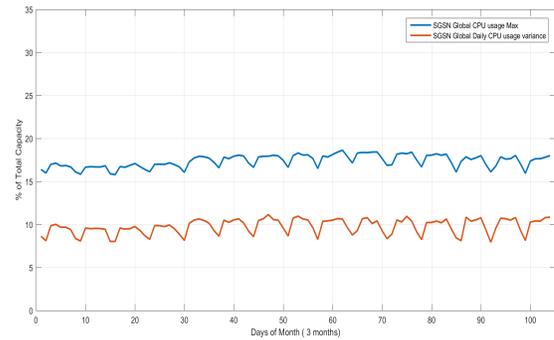
A. Measured Data

Load profiles were created for each node and on a global level for both networks to evaluate how load evolves over a period of time. Control Plane (CP) and Data Plane (DP) resources utilization were analyzed separately from these nodes. From Fig 2, it can be seen that the maximum CPU utilization on a daily scale, for both node level and global level, does not exceed 25 % of the total capacity. The main reason behind this is that nodes are optimized in such a way that in case of a failure, the load is taken by the other nodes in the same pool. It was also noticed that the CPU usage variance is almost constant with a slight decrease pattern on weekends for the whole measurement period. These maximum CPU usages along with daily CPU variance are good indicators to dimension nodes efficiently for failure cases. It was also observed that the control plane usage was higher for NW 1 while the data plane usage dominates for SGSN-MME in NW 2 and that is due to limited 4G services in NW 2 as it can be observed in Fig. 3. It was also observed that resource usage follow a repetitive pattern on daily scale. After analyzing both networks, we can draw the conclusion that days of month are pretty much the same on the basis of resource utilization, variation occurring primarily as a function of the time of the day. However, resource usage patterns differ from a network to another but repeat themselves on a daily basis within the same network as shown in Fig. 3, where each line indicates one day of a month.

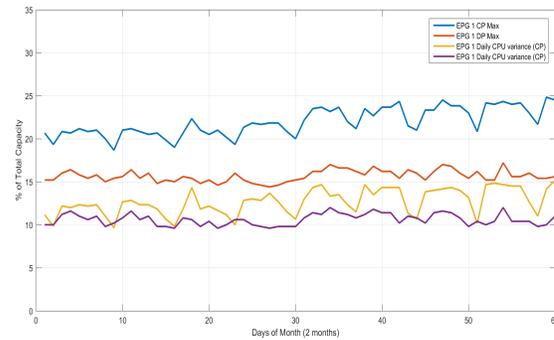
Fig. 4 plots the total CPU usage of SGSN-MMEs at global level for both studied networks. From the figure, it can be



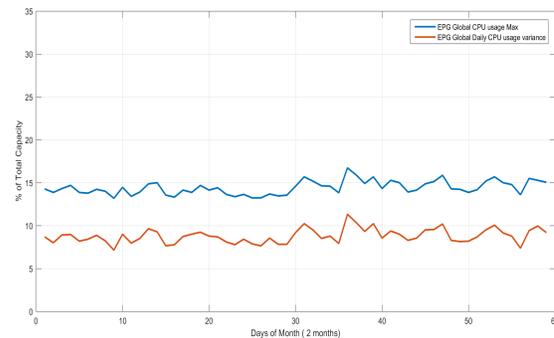
(a) SGSN-MME CP & DP maximum CPU usage and daily variance (NW 1).



(b) SGSN-MME Global maximum Total CPU usage and Daily variance (NW 1).

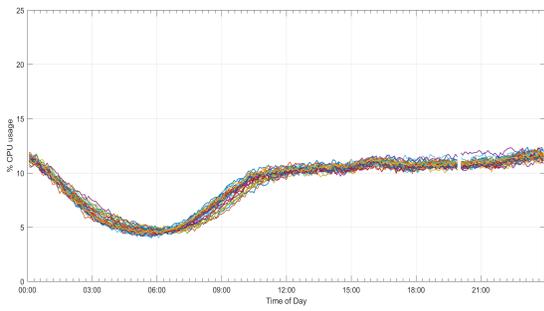


(c) EPG CP & DP maximum CPU usage and daily variance (NW 1).

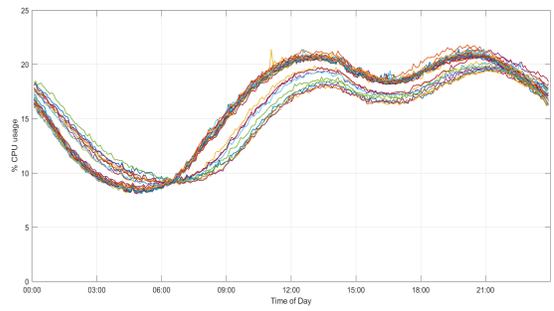


(d) EPG Global maximum Total CPU usage and daily variance (NW 1).

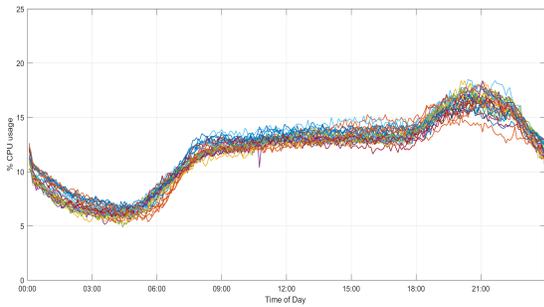
Fig. 2: Maximum and total CPU usages and daily variances for NW 1 Nodes.



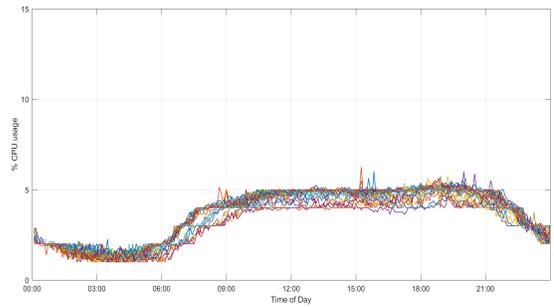
(a) SGSN-MME DP CPU Usage (NW 1)



(b) SGSN-MME CP CPU Usage (NW 1)

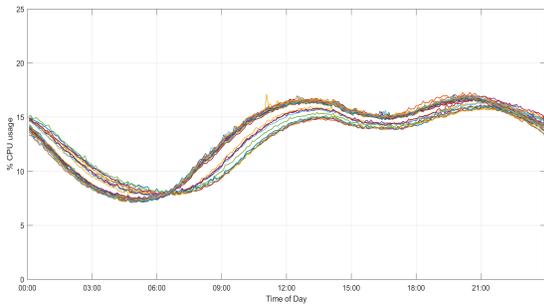


(c) SGSN-MME DP CPU Usage (NW 2)

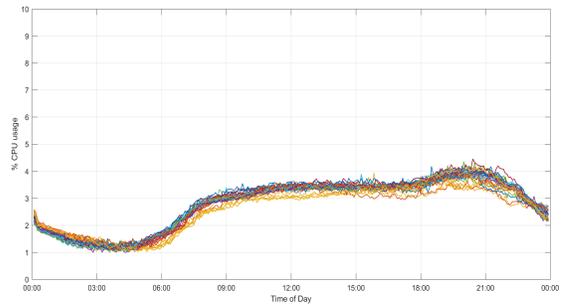


(d) SGSN-MME CP CPU Usage (NW 2)

Fig. 3: DP and CP CPU usage of SGSN-MME for both studied networks.

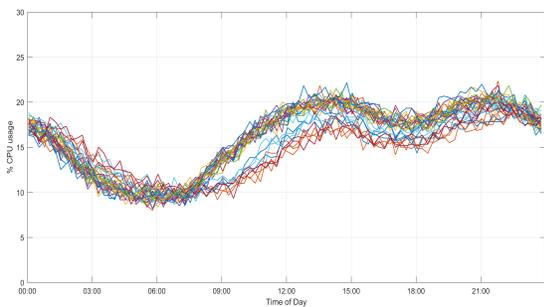


(a) SGSN-MME Global CPU usage (NW 1)

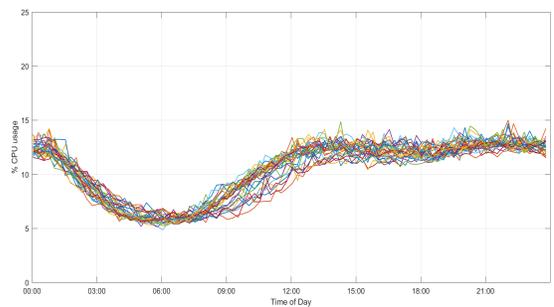


(b) SGSN-MME Global CPU usage (NW 2)

Fig. 4: SGSN-MME Global total CPU usage for both studied networks.



(a) EPG CP CPU usage on node level (NW 1)



(b) EPG Global CP CPU usage (NW 1)

Fig. 5: EPG CP CPU usage on node and global level for NW 1.

noticed that in NW 2 the utilization is lower compared to NW 1. It is because that NW 2 is upgraded lately for 4G services and most of the control plane processors are not used at all. In NW 2, only one EPG is utilized and its utilization level is quite low, making it impossible to draw any observations from the respective graphs and they are thus not included in this article. However for NW 1, EPG CP CPU usage on node and global level can be observed from Fig. 5.

From the studied networks, it is worth underlining the repetitive nature of the load as well as the fact that the best predictor of how load is likely to evolve is the time of day and the day of the week (e.g., weekend days) giving a further refinement opportunity. These observations are important for designing suitable scaling algorithms for networks deployed in virtual environments [6][13].

B. Simulation Results

Based on real load data, we evaluated how much actual capacity (i.e., number of blades) would be required, if the same load would have been supported in a virtualized environment. Three assumptions were made:

- 1) We can scale the allocated capacity according to actual load. In reality, one should assume a maximum 80 % load level (an industry-wide best practice) and take into account the granularity of allocated VMs.
- 2) There is a constant overhead due to virtualization both for control and data plane traffic, modeled as the factor α , with a value of 1.1 based on authors' prior empirical experience.
- 3) There is an additional overhead for data plane processing, which we modeled as β , to quantify the impact on throughput and latency.

Same kinds of calculation were carried out both for node level and global level. It can be seen that even with high β values the amount of generic hardware needed is still significantly below its native counterpart as shown in Fig. 6. NW 2 plots are not included in this article as the number of dedicated resources are low and horizontal scaling, the focus of this study, does not yield interesting results.

A number of interesting conclusions can be drawn from these simulations. First, the overhead, introduced by virtualization, plays a smaller role than expected; even with a supposedly low performance implementation, NFV will outperform native installations. This is due to a number of factors: virtualization reduces the need for dedicated resources to achieve geographical redundancy (i.e., nodes can be delivered at just two different locations datacenters to achieve geographical redundancy instead of three or more as today) and dynamically scale them according to load conditions; at the same time resource utilization can be kept at the optimal 80 % (which is a typical industry best practice level) most of the time, hence approximating optimal resource utilization levels.

Second, the most important technology to unlock the benefit of virtualization is optimized scaling, enabling allocation of resources on a need basis. An algorithm that can accurately

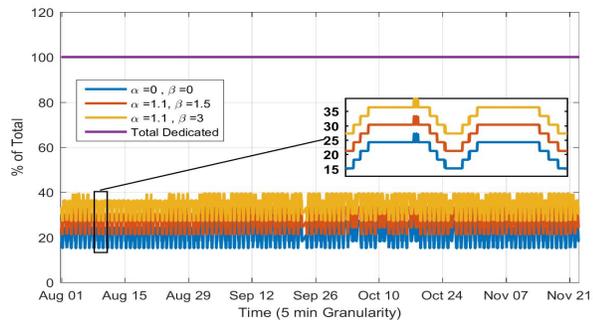
predict how load will evolve will hence be an essential ingredient of any commercial NFV deployment [12][13]. In this study, we assumed that horizontal scaling or scale out (i.e., adding or removing VMs based on load) is used. However, for NW 2 vertical scaling (i.e., resizing VMs based on load) would be more interesting and cost effective [13]. In general, horizontal scaling is more effective when amount of dedicated resources are high.

V. CONCLUSIONS

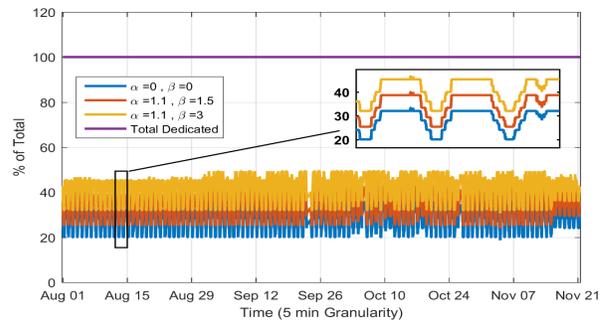
Based on our analysis of data from real-life mobile networks, the overall load is a fraction of installed capacity which applies both to global (all geo- redundant) and node level. Furthermore, resource utilization is different from a network to another but is highly correlated with time of day and follows the same pattern within the same network. It is concluded that virtualization with dynamic scaling of node size based on load is more cost efficient even if large virtualization overhead for data-plane is taken into account. This defines an interesting research area where future research shall focus. Additionally, for networks with small numbers of dedicated resources, vertical scaling (resizing the VMs) will be more cost-effective than horizontal scaling (adding more VMs) [13]. Future research work includes identifying any other possible correlation related to load to develop dynamic scaling algorithm for resources based on utilization and verify it on real network conditions.

REFERENCES

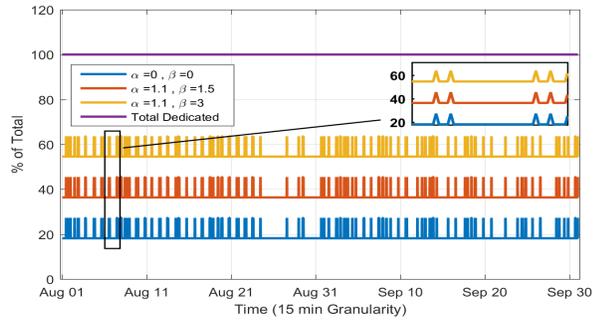
- [1] ETSI, "Network Functions Virtualisation: An Introduction, Benefits, Enablers, Challenges & Call for Action," NFV - Introductory white paper, Oct. 2012
- [2] K. Tan et al., "Sora: High performance software radio using general purpose multi-core processors," in Proc. USENIX Int. Symp. NSDI, 2009, pp. 7590.
- [3] R. Shah, "Network I/O Virtualization : Challenges and Solution Approaches," Seminar Report Indian Institute of Technology, Apr. 2014.
- [4] B. Johnson and R. Deshmukh, "Boosting I/O performance for virtualized servers," Dell Power Solutions, No. 3, 2012.
- [5] Intel, "PCI-SIG SR-IOV Primer An Introduction to SR-IOV Technology," Revision 2.5, Jan. 2011.
- [6] T. Taleb, M. Corici, C. Parada, A. Jamakovic, S. Ruffino, G. Karagiannis, and T. Magedanz, "EASE: EPC as a Service to Ease Mobile Core Network," in IEEE Network Magazine, Vol. 29, No. 2, Mar. 2015. pp.78-88.
- [7] T. Taleb, "Towards Carrier Cloud: Potential, Challenges, & Solutions," in IEEE Wireless Communications Magazine, Vol. 21, No. 3, Jun. 2014. pp. 80-91.
- [8] T. Taleb, A. Ksentini, M. Chen, and R. Jantti "NFV-based Dynamic Service Chaining for Emerging Mobile Social Media Applications", to appear in IEEE Trans. on Wireless Communications.
- [9] T. Taleb, A. Ksentini, and R. Jantti, "Anything as a Service for 5G Mobile Systems" to appear in IEEE Network Magazine.
- [10] T. Taleb, A. Ksentini, and B. Sericola, "On Service Resilience in Cloud-Native 5G Mobile Systems", to appear in IEEE J. Select. Areas in Communications.
- [11] F.Z. Yousaf, P. Loreiro, F. Zdarsky, T. Taleb, and M. Leibs, "Cost Analysis of initial deployment strategies of a Virtual Network Infrastructure in a Datacenter," in IEEE Communications Magazine, Vol. 53, No. 12, Dec. 2015, pp. 60 - 66.
- [12] F.Z. Yousaf and T. Taleb, "Fine Granular Resource-Aware Virtual Network Function Management for 5G Carrier Cloud," to appear in IEEE Network Magazine.



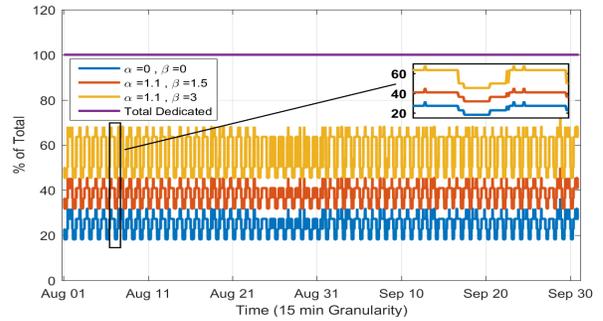
(a) SGSN-MME (Node Level)



(b) SGSN-MME (Global Level)



(c) EPG (Node Level)



(d) EPG (Global Level)

Fig. 6: Simulation Results

- [13] S. Dutta, T. Taleb, and A. Ksentini, "QoE-aware Elasticity Support in Cloud-Native 5G Systems," in IEEE ICC '16, Kuala Lumpur, Malaysia, May 2016.