# Performance Analysis of SDN and NFV enabled Mobile Cloud Computing

Joseph Billingsley*, Wang Miao*, Ke Li* Geyong Min*, and Nektarios Georgalas†

*Department of Computer Science, University of Exeter, UK

Email: {jb931, wang.miao, k.li, g.min}@exeter.ac.uk

†Research and Innovation, British Telecom, UK

Email: nektarios.georgalas@bt.com

*Abstract*—Mobile Cloud Computing (MCC) is regarded as a promising method to extend the battery life, increase the data storage and enhance the processing power of mobile devices. The technologies, e.g. Software Defined Networking (SDN) and Network Function Virtualisation (NFV) have been deployed in MEC to simplify the network management and accelerate mobile service deployment. There have been some interesting research findings appeared in the literature regarding the performance of SDN and NFV in MCC, however most of the existing work only considers these technologies in isolation and pays little attention to their cooperative and complementary relations in practical deployments. In order to achieve a deeper understanding of future MCC, a comprehensive analytical model is developed in this work to investigate the performance of MCC in the presence of both NFV service chains and SDN networks. The proposed model is capable of investigating the interactions between SDN and NFV when they share the same underlying physical infrastructure. The end-to-end latency is derived based on the developed model with different scales of service deployments and network configurations. Comprehensive simulation experiments are conducted and the results demonstrate that the proposed analytical model matches well with the simulation experiments. In addition, the proposed analytical model is used as a useful tool to investigate the impact of the centralised SDN control on the performance of the NFV traffic transmission.

## I. INTRODUCTION

Emerging mobile services such as Virtual Reality, 4K video, and tactile internet, consume incredible amounts of compute, storage and bandwidth resources [?]. However, due to the inherent constraints of their size, weight and power, mobile devices become struggle to meet the strict resource requirements of new applications. Mobile Cloud Computing (MCC) [?] [?] has been considered as a key technology for mobile devices to address this issue. Through migrating the local applications and service to a mobile cloud datacenter, MCC brings mobile devices the benefits of extending the battery life, increasing the data storage and enhancing the processing power. To realise this ambition, Software Defined Networking (SDN) and Network Function Virtualisation (NFV) have been regarded as two promising and complementary technologies in MCC datacenter to simplify the datacenter network management and improve the resource utilisation and service flexibility.

SDN is a new networking architecture that can simplify the network management and accelerate network innovation. It is implemented by decoupling the network control from the underlying network infrastructure and creating a software-programmable controller. A logically centralised SDN controller maintains a global view of the network, helping the network operator to design network services and determines how packets should be routed through the network [?] [?]. With the centralised network control, networks intelligence is migrated from the underlying network devices to SDN controller, enabling network operators to manage the entire network consistently and holistically.

NFV is a novel network architecture which allows for flexibility in service provisioning. Traditionally, services are constructed by connecting chains of purpose built computers each performing a particular function. These may be traditional data centre functions such as firewalls and load balancers, or mobile communications functions such as the Packet and Service gateways in the 4G Evolved Packet Core. NFV decouples these functions from the hardware by implementing these network functions in software on virtual machines (VMs). These Virtual Network Functions (VNFs) can be moved, scaled or destroyed on demand, allowing for efficient placement and allocation of resources, significantly accelerating the deployment of new services.

SDN and NFV are often considered complementary technologies in practical deployments [?]. For instance, when a new MCC service needs to be deployed in the cloud datacenter, the cloud management system firstly design a service chain and leverages Virtual Network Function (VNF) manager to deploy VNFs in the underlying Virtual Machines (VMs) or containers. After initiating the VNFs, the address of the VMs or containers as well as the NFV chain will be sent to SDN controller, which is responsible for establishing the connections between difference VNFs and collecting the stochastic information cloud management system for service optimisation. From the above example, it can be seen that SDN plays an important role in the deployment, management and optimisation of overall lifecycle of NFV service provisioning. Therefore, it is necessary to jointly consider SDN and NFV in the datacenter network management and service provisioning. For the system optimisation, analytical models can provide insight of system operation by formally defining the interactions among key parameters such as the scale and resources utilisation of the datacenter, the traffic generated by the end devices and the Quality of Service (QoS) performance that should be satisfied. There have been some

research efforts to analyse the performance of SDN and NFV network architecture [?] [?] [?] [?] [?]. For instance, for modelling the performance of SDN networks, Longo et al. [?] proposed a model to investigate the reliability performance of SDN network, which is based on a two layer network management architecture. Azodolmolky et al. [?] exploited network calculus to investigate the worse case delay performance of SDN network as well as minimum buffer size required for a given delay constraint. To capture the burstness of the network traffic, Wang et al. [?] developed an analytical model to investigate an SDN architecture with the traffic following the Markov-modulated Poisson Process (MMPP). A high and low priority queue model was proposed in this work to capture the features of the practical SDN deployment. For analysing the performance of NFV architecture, Prados-Garzon et al. [?] designed an analytical model to investigate the average response time of a single NFV service provisioning. Gebert et al. [?] modeled each step in the packet processing pipeline of a Linux x86 system, which is based on single VNF deployment scenario. To analyse the stochastic performance of multiple VNF scenario, an analytical model was proposed in [?] to exploit the stochastic network calculus to investigate of the end-to-end performance of an NFV service provisioning, which could obtain the worse case of network transmission for a given QoS requirements. Although, these existing works provide some insights into the performance of SDN and NFV network architecture in various network scenarios, they seldom jointly consider these two technologies in the performance analysis. From the perspective of service deployment and provisioning, SDN and NFV are complementary technologies and always deployed together. Therefore, it is important to investigate the performance of network infrastructure with both SDN and NFV support, especially identifying how their interactions can affect the performance of service provisioning. To the best of our knowledge, only Fahmin et al. [?] have considered both NFV and SDN in their analytical model. However the network infrastructure adopted in [?] consists of only one switch and one VNF, which can be hardly to be applicable to a large scale datacenter network. In order to reap the benefits of SDN and NFV for MCC applications, there is an urgent need to develop a novel analytical model which can jointly consider two complementary technologies in a large-scale datacenter network.

To fill this gap, a comprehensive analytical model is proposed in this work to investigate the performance of SDN and NFV enabled MCC datacenter networks. To capture the unique features of real-world SDN and NFV deployments a network architecture is firstly abstracted in this study with multiple NFV chains and a virtualised SDN implementation, where the SDN controller determines how traffic is routed among the VNFs. The analytical model is developed with the aim of understanding the interactions between SDN and NFV when they are deployed under the same underlying physical infrastructure, e.g the impact of the length of NFV service chain on the traffic engineering performance of SDN networks. The end-to-end performance in terms of average latency is ob-
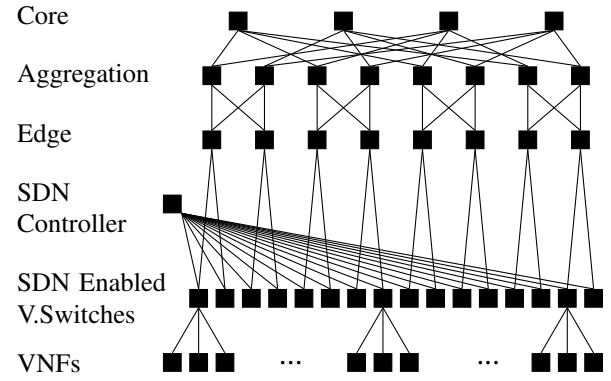


Fig. 1. An example SDN and NFV enabled fat-tree network with 4 ports for each hardware switch and 3 for the virtual switches.

tained by the developed model and validated through extensive simulation experiments under different network configurations. In addition, the proposed analytical model could be used as an effective tool to optimise the design of services and networks in MCC datacenters.

The remainder of this paper is organised as follows. Section II discusses the details of the network architecture that is modelled in this work. In Section III we derive the analytical model for the network. Section IV validates the accuracy of this model with extensive simulation experiments. Finally Section V concludes the paper.

## II. Network Architecture

Based on the general working mechanism of SDN and NFV technologies [?], a network architecture is abstracted in this study as shown in Fig. 1, where multiple NFV chains are deployed and linked by a virtualised SDN implementation. With NFV deployment, a MCC service is provided in the form of a service chain which is formed by several VNFs. The packets of MCC devices pass through each of the VNFs in sequence. In the abstracted network architecture, multiple NFV chains can coexist and different NFV service chains have different numbers of VNFs, each of which has different packet processing capability.

When a MCC service is launched in the cloud datacenter, the cloud management orchestrator or manager analyses the performance and functional requirements, and initiate VNFs in the underlying virtualised server. Once the VNFs are deployed, the cloud management orchestrator coupled with SDN controller link the individual VNFs to form an NFV service chain to provide MCC service for end devices. During this process, SDN controller is responsible for configuring the routing table of the underlying network switches, collecting the related service parameters, and sending the collected data to cloud management orchestrator for service optimisation. Each SDN enabled switch has a flow table containing instructions on how to route the received packets. Due to the physical storage limitation, it is impossible to store the instructions for all possible destinations. Therefore, once the local switch receives a packet that doesnot match routing table, a request

will be sent from switch to the controller to consult on how to process the received package. After a series of computation, SDN controller responses this request for further process. Within the abstracted network architecture, we consider an SDN architecture where only the virtual switches connect to the SDN controller. This architecture is representative of those used in industry, most notably a comparable architecture is used in VMWare's NSX solution [?].

According to [?], the network topology that supports the communication among different architecture components, e.g. VNFs and SND controller, is mainly based on a fat-tree structure, which is formed by three layers of switches, core, aggregation and edge switches. In the most of the modern datacenter, the switches at the edge layer are connected to Top-of-Rack (ToR) switches, and the VNFs are hosted in the VMs or containers of datacenter servers. As shown in Fig. 1, the fat-tree topology is defined by the number of ports at each switch. Let $k$ denote the number of ports for each physical switch and $k_v$ be the number of ports for each virtual switch. Each core switch connects to the switch of the pod, which contains two layers of switches, aggregation and edge respectively. Within the pod, edge switches are fully connected to aggregation switches. In addition, each edge switch is connected to $k/2$ servers. Each server contains a virtual switch connected to $k_v$ VNFs. Based on the connection relationship between core, aggregation, edge, virtual switches, and VNFs, a three layer $k$ pot fat-tree topology has $(k/2)^2$ core switches, $k$ pod, $k^2/2$ aggregation switches, $k^2/2$ edge switches, $k^2$ virtual switches, and $(k^3/4) \cdot k_v$ VNFs.

Following the work of SDN and NFV performance analysis in [?] [?] [?], it is assumed that packets coming from different MCC users will be independent from each other and service time of processing a packet is also independent on earlier packets. Hence the arrival and service rates of packets in this study follow independent probability distributions. For each NFV chain, the traffic entering the MCC datacenter follows an independent Poisson process with a mean rate of $\lambda$ packets per second. Each physical/virtual node, e.g. switch, VNF and the controller, provide the services for the coming packets according to an independent Poisson process with service rates of $\mu_s$, $\mu_v$ and $\mu_c$ packets per second respectively. If a packet fails to match the routing table in the SDN-enabled virtual switches, the information of this packet will be forward to SDN controller for further processing. Let $p_m$ denote the probability of there is no routing entry for the incoming packet.

## III. ANALYTICAL MODEL DERIVATION

In this section, we will present the methodology and approaches to derive the analytical model of SND and NFV enabled MCC datacenter. The model to be designed will be capable of analysing the end-to-end service performance with the coexistence of multiple NFVs, each of which will have different number of VNFs with respect to different MCC services. The impact of the centralised control of SDN controller on the end-to-end performance provisioning is also studied in the developed model. With the aim of increasing the readability of model derivation, we firstly consider the case of single NFV deployment scenario in the analytical derivation, then the simplified version will be extended to the scenario of multiple NFV service chains.

*1) Single NFV deployment Case:* According to the widely used Equal-Cost Multi-path Routing (ECMR) [?] in datacenter network, the end-to-end latency is dependant on the probability that a packet will visit a certain layer of switches and the waiting times in each layer. The end-to-end latency can be written as,

$$L_t = \prod_{i=0}^{3} L_i(\lambda_i, \mu_i) p_i \tag{1}$$

where "$i = 0$" represents the virtual switch layer, and "$i = 1, 2, 3$" denotes the edge, aggregation and core layers respectively. $L_i(\lambda_i, \mu_i)$ denotes the end to end latency when the packets need to travel to the $i$th layer network switches. Similarly $p_i$ represent the probability that a packet could reach the $i$th layer switches. $L_i(\lambda_i, \mu_i)$ is the sum of the latency the packets experience from VNF nodes to the $i$th layer switches, which can be calculated by,

$$L_i(\lambda_i, \mu_i) = w_f + \sum_{j=0}^{i} w_j \tag{2}$$

where $w_f$ is the processing latency within a VFN. $w_j$ is the latency at the $j$th layer of the switches. For the virtual switches, the latency, $w_0$, has two parts, the latency at the virtual switch and latency at the SDN controller. Therefore, $w_0 = p_c w_c + (1 - p_c) w_v$, where $p_c$ is the probability that a packet will be forward to SDN controller for the routing decision making. If the source and destination VNFs share the same virtual switch, the packets between two VNF will not visit a higher layer switch. Let the probability of a packet only visiting a virtual switch, $p_0$, denote is the probability that the source and destination VNFs are under the same virtual switch. $p_0$ can be calculated by

$$p_0 = \frac{k_f - 1}{n_v - 1} \tag{3}$$

where $n_v$ denotes the total number of VMs in the datacenter.

Based on the topology of fat-tree structure, the higher layer a packet can reach means that the more destinations this packet could visit. If the switches that a packet visits is at edge layer, then the probability that the destination VNF is under the same edge switch can be calculated by excluding destinations that could be visited via a shorter path. In this case, the short path would be virtual switches. Therefore, $p_1$ can be derived from the following equation,

$$p_1 = \frac{(k/2 - 1) \cdot k_f}{n_v - 1} \tag{4}$$

Following the method of deriving $p_1$, the probability of visiting the aggregate and core layers can be calculated by,

$$p_2 = \frac{(k/2 - 1) \cdot k \cdot k_f}{2 \cdot (n_v - 1)} \tag{5}$$

$$p_3 = \frac{n_v - (k/2)^2 \cdot k_f}{n_v - 1} \tag{6}$$

At the virtual switch, the SDN controller will only be consulted if the destination VNF is located in another physical server, and the packet does not match routing table of virtual switch in that server. Then, the probability that the packets will be sent to SDN controller for routing computation can be calculated by,

$$p_c = (1 - p_0) \cdot p_m \tag{7}$$

After obtaining the probability that a packet will be processed at the different layers of switches and SDN controller. The following subsection derives the waiting time at each component of the routing path. According to [**?**], the waiting time for a M/M/1 queue is obtained by

$$w(\mu, \lambda) = \frac{1}{\mu - \lambda} \tag{8}$$

where $\mu$ is the service rate and $\lambda$ is the arrival rate for a M/M/1 queue. In the following, we aim to calculate the arrive rate at the VNFs, SDN controller and different layers of switches. As the destination VNFs are evenly distributed over the VNFs, each destination VNF will receive an equal proportion of packets from other VNF. Hence the traffic arrives at the destination VNF at the rate of $\lambda_f$. Virtual switches realise the communications among VMs, so virtual switches can receive packets from three sources: 1) packets generated by VNFs on the server that the virtual switch locates; 2) the traffic generated by the VNFs in the other servers; and 3) the packets feed back from the SDN controller. Then the arrival rate at the virtual switch can be calculated by,

$$\lambda_0 = \lambda_f (1 + \frac{n_v - k_v}{n_v - 1} + p_c) \tag{9}$$

The arrival rate for the edge switches can be achieved similar to the virtual switch. It should be noticed that packets that are intended for destinations on the same server do not need to visit the edge switch. After a series of derivation, the arrival rate at the edge switch can be calculated by,

$$\lambda_1 = \frac{\lambda_f \cdot k \cdot k_v}{2(n_v - 1)}(n_v - k_v + \frac{(2n_v - k) \cdot k_v}{2}) \tag{10}$$

According to the MCMR routing protocol in the datacenter network, the traffic will be balanced among aggregate switch in a pod. And the VNFs sharing the same virtual or edge switches will not reach the aggregation switches. Then arrival rate at the aggregate switch can be computed by,

$$\lambda_2 = \frac{\lambda_f \cdot k \cdot k_v}{2(n_v - 1)}\left(2n_v - \frac{k}{2}\left(k_v - \frac{k \cdot k_v}{2}\right)\right) \tag{11}$$

As all VNFs are connected by the core switches, the arrival rate at each core switch is the portion of traffic that their destination VNFs cannot be reached by edge or aggregation switches. Based on MCMR protocol, the traffic leaving the aggregation layer will be evenly split among different core switches. Therefore the arrival rate at the core switch is calculated by,

$$\lambda_3 = \frac{\lambda_f \cdot p_o \cdot n_v}{(k/2)^2} \tag{12}$$

Finally, let us calculate the traffic rate for the SDN controller. It can be observed that the packets visiting the SDN controller do not change the arrival rate of higher level switches, e.g. edge or aggregation switches; and a portion of the traffic in the virtual switch will be sent to SDN controller for routing decision making. Given the number of the virtual switch ($n_v = k^2$), then the arrival rate at the SDN controller is computed by,

$$\lambda_c = \lambda_f \cdot p_c \cdot n_v \tag{13}$$

By substituting the arrival rates (Eqs. 10 to 12) and service rates of each network component into M/M/1 latency equation, we can obtain the average waiting time at VNF, and different layers of switches. To calculate the latency caused by the round-trip between virtual switches and SDN controller, two latency components should be considered, the latency in the SDN controller and the latency in the virtual switch when the packet return back from the SDN controller. So, the additional latency caused by the SDN controller could be calculated by $w_d = w_c + w_v$. Through taking the probabilities of the different paths and the mean waiting time into Eq. (1), we could obtain the end to end latency for the single VNF deployment scenario.

*2) Multiple NFV Deployment Case:* Although there are several researches investigating the performance of NFV, existing research only considered the case of a single NFV service deployment, paying little attention to the multiple NFV deployment. Given the fact that datacenter infrastructure is always used to simultaneously support multiple services, it is necessary to investigate the performance of SND and NFV enabled datacenter network with the different scale of service deployment. In this subsection, we will extend the simplified NFV deployment scenario into multiple NFV case.

Let $N_s$ denote the number of NFV services deployed in datacenter and $K_i$ represent the length of $i$th service. Along the NFV chain, the packets will visit each VNF in sequence before arriving at the end VNF. So the VFN receives the traffic forwarded from the previous VNF, and produces and send the new packets to the next VNF. Therefore, each VNF except the last one would receive and generate packets simultaneously. From the perspective of the network transmission, the network infrastructure will receive the packets from $K_i - 1$ VNFs for the $i$th service and the effective network traffic rate that $i$th service chain generates can be written as $\lambda_{i,f}^e = \lambda_{i,f} \cdot (K_i - 1)$, where $\lambda_{i,f}$ is the traffic rate of the $i$th service and the $\lambda_{i,f}^e$ is the effective traffic rate that the network infrastructure receives. In the abstracted analytical model, different network service would have different length of the service chain. In

addition, the end-to-end latency a packet persisted in the network is dependent on both the type and service chain length of the network services. As NFV services are independent with each other, let $p_{i,s}$ denote the probability that a packet network device receives is from the $i$th service. The effective network traffic the network infrastructure receives could be obtained by

$$\lambda_f^e = \sum_{i=1}^{N_s} p_{i,s} \cdot \lambda_{i,f} \cdot (K_i - 1) \tag{14}$$

Similar to the single NFV deployment case, the end-to-end latency will be the probabilistic sum of the time spent taking each path. By inserting the effective traffic rate in Eq. (14) in to Eqs. (9-13), we could calculate the effective network traffic in each network layer. Given the service rates $\mu_v$, $\mu_e$, $\mu_a$, and $\mu_s$, the latency, $w_j$ could be obtained for different network layer. After calculating the probability that a packet visits a certain network layer, the end-to-end latency can be calculated by Eqs. (??) and (1). For convenience, pseudocode for the entire process is given in Algorithm 1.

---

**Algorithm 1** Calculation of Average Latency of SND and NFV-enabled MCC Datacenter Networks

---

1: Calculate the effective network traffic: $\lambda_f^e$     (*Eq. (14)*)
2: Calculate the traffic rates: $\lambda_i$     (*Eqs. (??-??)*)
3: Calculate the probability: $p_i$     (*Eqs. (9-13)*)
4: Calculate the waiting time: $w_j$ and $w_f$     (*Eq. (8)*)
5: Calculate the latency for the $i$th path: $L_i$     (*Eq. (2)*)
6: Calculate the latency for the end to end transmission: $L_t$ (*Eq. (1)*)

---

## IV. VALIDATION

To verify the accuracy of the analytical model, a discrete event simulator has been built using OMNeT++ [?] to simulate a NFV and SDN enabled datacentre network. Each simulation experiment was run until the network reaches its steady state where further network cycles do not change the collected statistics appreciably. Comprehensive simulation experiments were conducted to validate the performance of the proposed analytical model under different network configurations. However for the sake of specific illustration only a selection of tests are presented here and the results comparison between the analytical model and simulation experiments are presented in terms of the average end-to-end latency.

In practice a datacentre can contain on the order of tens of thousands of servers [?], with each switch supporting 1 to 100Gbits/s traffic a second. It is hardly to simulate the scale of datacenter network in lab environment. Therefore, a scaled down version of a typical datacentre is modelled with the following parameters,

- $k = \{4, 6, 8\}$, $k_v = 2$ and $p_m = \{0.1, 0.5, 1\}$
- The service rate of the switches and SDN controller are set to be 40 packets per second ($\mu_v = 40$, $\mu_c = 40$)
- The service rate of the VNFs is set to be 20 packets per second ($\mu_f = 20$)

- Services are selected with equal probability
- Case I: The network holds one service with two VNFs
- Case II: The network holds multiple services ($N_s = \{1, 2, 4\}$ and $K_i = \{3, 5, 7, 9\}$)

Figs 2 to 5 depict mean message latency predicted by the model plotted against those provided by simulation experiments for a range of parameter settings. For the model, results are only shown where the network is in a steady state, i.e. where the arrival rate is lower than or equal to the service rate for all queues. The figures demonstrate that the simulation results closely match those predicted by the model. The tractability and accuracy of the analytical model make it suitable for analysis of next generation NFV and SDN enabled Mobile Cloud computing datacentres.

## V. CONCLUSION

In this paper we have presented a comprehensive analytical model to investigate the performance of a SDN and NFV enabled MCC datacenter. Firstly, we abstracted a network architecture to capture the behaviour of SDN and NFV when they are implemented in a datacenter. Based on the abstracted network architecture, a novel analytical model was developed with the aim of obtaining the end-to-end latency performance for service provisioning. The proposed model is capable of investigating the interactions between SDN and NFV when they share the same underlying physical infrastructure. Extensive simulations were conducted to validate the performance of the proposed analytical model. Simulation results demonstrated that the proposed analytical model matches well with simulation experiments. The proposed resulting analytical model is fast and accurate, and hence suitable for determining the optimal design of services and networks in large scale Mobile Cloud computing datacenters.

## REFERENCES

[1] J. G. Andrews, S. Buzzi, W. Choi, S. V. Hanly, A. E. Lozano, A. C. K. Soong, and J. C. Zhang, "What will 5g be?" *IEEE Journal on Selected Areas in Communications*, 2014.

[2] H. Kim and N. Feamster, "Improving network management with software defined networking," *IEEE Communications Magazine*, 2013.

[3] S. Hares and R. White, "Software-defined networks and the interface to the routing system (I2RS)," *IEEE Internet Computing*, 2013.

[4] F. Longo, S. Distefano, D. Bruneo, and M. Scarpa, "Dependability modeling of software defined networking," *Computer Networks*, 2015.

[5] W. Miao, G. Min, Y. Wu, H. Wang, and J. Hu, "Performance modelling and analysis of software-defined networking under bursty multimedia traffic," *TOMCCAP*, 2016.

[6] J. Prados-Garzon, P. Ameigeiras, J. J. Ramos-Muñoz, P. Andres-Maldonado, and J. M. López-Soler, "Analytical modeling for virtualized network functions," in *ICC'17: IEEE International Conference on Communications*, 2017.

[7] S. Gebert, T. Zinner, S. Lange, C. Schwartz, and P. Tran-Gia, "Performance modeling of softwarized network functions using discrete-time analysis," in *ITC'16: International Teletraffic Congress*, 2016.

[8] S. Azodolmolky, P. Wieder, and R. Yahyapour, "Performance evaluation of a scalable software-defined networking deployment," in *EWSDN'13: European Workshop on Software Defined Networks*, 2013.

[9] A. Fahmin, Y. Lai, M. S. Hossain, Y. Lin, and D. Saha, "Performance modeling of SDN with NFV under or aside the controller," in *Fi-Cloud'17: International Conference on Future Internet of Things and Cloud Workshops*, 2017.
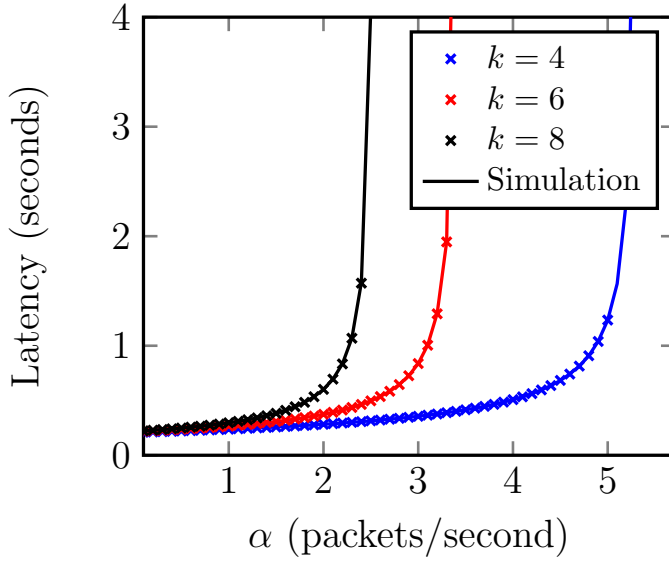
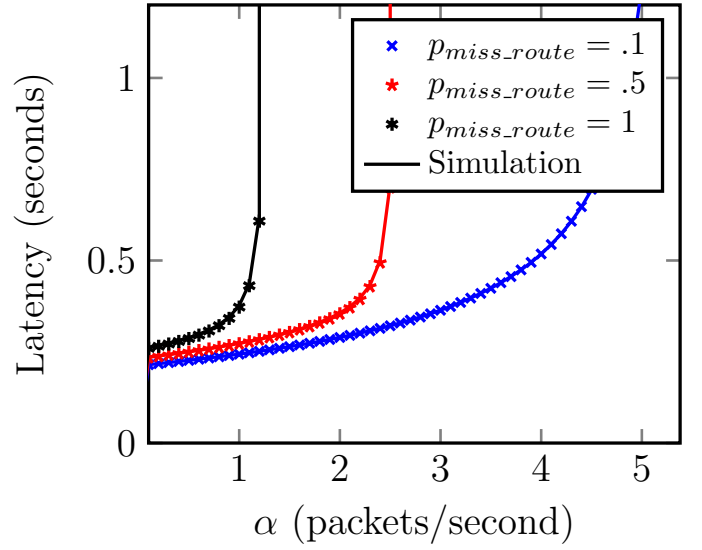Fig. 2. Latency predicted by the model and simulation for different numbers of ports ($k$).



Fig. 3. Latency predicted by the model and simulation with different proportions of packets visiting the SDN controller ($p_{miss\_route}$).
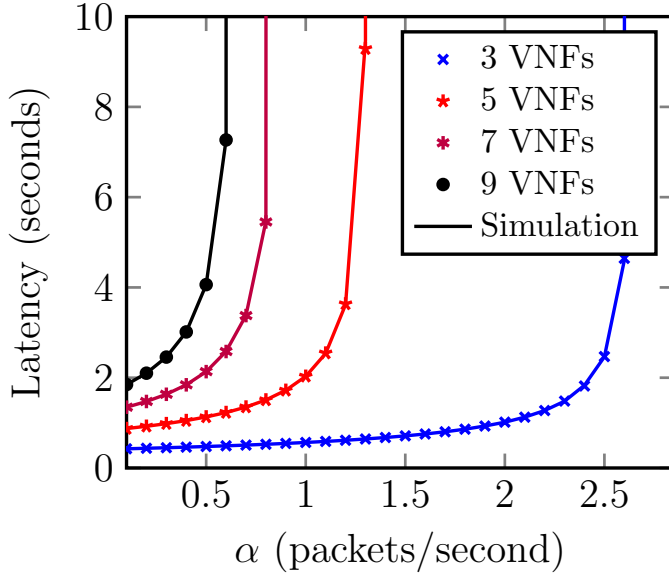


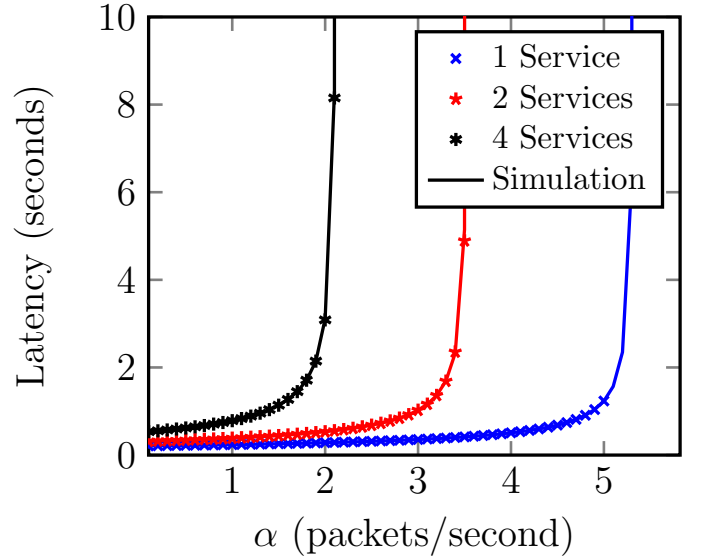Fig. 4. Latency predicted by the model and simulation for different length service chains.



Fig. 5. Latency predicted by the model and simulation for several services with different length service chains.

[10] VMware, "Network virtualisation and security platform - nsx," 2018, [Online; accessed 2018-05-04]. [Online]. Available: https://www.vmware.com/uk/products/nsx.html

[11] Cisco, "Cisco global cloud index: Forecast and methodology, 2016-2021," 2018, [Online; accessed 2018-05-04]. [Online]. Available: https://www.cisco.com/c/en/us/solutions/collateral/service-provider/global-cloud-index-gci/white-paper-c11-738085.html

[12] L. Kleinrock, *Theory, Volume 1, Queueing Systems*. Wiley-Interscience, 1975.

[13] A. Varga and R. Hornig, "An overview of the omnet++ simulation environment," in *SimuTools'08: International Conference on Simulation Tools and Techniques for Communications, Networks and Systems & Workshops, year = 2008,*.

[14] J. Hamilton, "Aws re:invent 2016," 2016, [Online; accessed 2018-05-04]. [Online]. Available: https://www.youtube.com/watch?v=AyOAjFNPAbA