

Improving Traffic Forecasting for 5G Core Network Scalability: A Machine Learning Approach

Imad Alawe, Adlen Ksentini, Yassine Hadjadj-Aoul, and Philippe Bertin

ABSTRACT

5G is expected to provide network connectivity to not only classical devices (i.e., tablets, smartphones, etc.) but also to the IoT, which will drastically increase the traffic load carried over the network. 5G will mainly rely on NFV and SDN to build flexible and on-demand instances of functional networking entities via VNFs. Indeed, 3GPP is devising a new architecture for the core network, which replaces point-to-point interfaces used in 3G and 4G by producer/consumer-based communication among 5G core network functions, facilitating deployment over a virtual infrastructure. One big advantage of using VNFs is the possibility of dynamic scaling, depending on traffic load (i.e., instantiate new resources to VNFs when the traffic load increases and reduce the number of resources when the traffic load decreases). In this article, we propose a novel mechanism to scale 5G core network resources by anticipating traffic load changes through forecasting via ML techniques. The traffic load forecast is achieved by using and training a neural network on a real dataset of traffic arrival in a mobile network. Two techniques were used and compared: RNN, more specifically LSTM; and DNN. Simulation results show that the forecast-based scalability mechanism outperforms the threshold-based solutions, in terms of latency to react to traffic change, and delay to have new resources ready to be used by the VNF to react to traffic increase.

INTRODUCTION

Network virtualization has paved the way for the fifth generation (5G) to meet high flexibility and agility expectations. Network virtualization allows running complex network functions (namely, virtual network functions — VNFs) on top of a virtualized infrastructure, hosted at central or edge telco cloud. Chained together, the VNFs allow building flexible and on-demand virtual networks. Such flexibility is highly required in 5G, where services need to be created on demand and for a given duration. Relying on network virtualization, standardization groups, like the European Telecommunications Standards Institute (ETSI) and Third Generation Partnership Project (3GPP), are multiplying efforts to build a new environment for 5G. ETSI, with the network functions virtualization (NFV) reference archi-

tecture [1], has created a reference model to orchestrate and manage VNFs running on top of virtual infrastructure. Meanwhile, 3GPP is developing specifications to deploy 3GPP elements, including the core network (CN) functions and some parts of the radio access network (RAN) functions [2, 3], on top of a virtualized infrastructure. Particularly, the envisioned 5G CN has been devised with the aim to run the new functions on a virtualized environment. One option of the 5G core architecture is to replace classical 3GPP interfaces (i.e., point-to-point) by a bus-like communication to interconnect the CN elements, following the concept of web services.

VNF dimensioning is one of the main challenges of running CN elements in a virtualized environment in terms of resources, which highly impacts performance [4]. Indeed, the virtualized CN elements should perform similarly to their hardware version, supporting the same quality of service (QoS): supported number of users, provided data rate, and so on. To reach the same performance as the hardware version, overprovisioning the allocated resources for a VNF could be envisioned, but it might be costly. Alternatively, adaptive provisioning of virtual resources could represent an interesting solution where the VNF is oversized (by increasing the number of CPUs or by multiplying the number of virtual machine, VM, instances of the VNF) when traffic load is growing and downsized when traffic is decreasing. Such a solution is known as a scale-up/down (if the number of CPU is increased/decreased) or scale-out/in processes (if the number of VMs serving the same VNF is augmented/reduced), and is often used in cloud computing. To scale VNF resources, a solution to make scaling decisions is needed. While in cloud computing the decision is mainly based on system-level information, like consumed CPU or memory, obtained via the virtualization platform, mobile networks require service-level information, like the number of users connected to the system or the traffic load, which indicate if the VNF is overloaded or not.

In this article, we are interested in the scalability of a 5G CN key element, namely the access and mobility management function (AMF). The latter represents the bottleneck of the mobile network control plane, as it handles user equipments' (UEs') attach requests. Accordingly, VMs hosting

the AMF should be well dimensioned in order to avoid increasing the attach duration of UEs and the percentage of rejected requests. The AMF should be oversized (scale-out) if the number of attach requests is high and downsized (scale-in) if the number of such requests is low. In [5], we proposed a mechanism to scale the AMF using a threshold-based algorithm, which was derived using control theory. The proposed mechanism uses a threshold (the number of attach requests received by the AMF) to make the decision to scale-out or -in. However, due to the time needed to scale-out (i.e., the time needed for a new VM to be operational), the proposed solution suffers from delay in reacting to an increase in the number of attach requests, which leads to an increase of the attach procedure duration, and hence the rejection of some requests. To overcome this issue, we propose in this article to use machine learning (ML) to forecast requests' arrivals, and hence anticipate the scale-out process, avoiding the rejection of requests and keeping the attach duration low.

The remainder of this article is organized as follows. The following section introduces the necessary background to understand the article, focusing on the 5G CN architecture and related works. Then we first introduce our motivations to ensure AMF scalability; next, the ML model and its usage to predict the user attach requests, and the proposed ML-based mechanism to scale-out and -in the AMF. We then present the result of using the ML algorithm to predict the user attach based on the Telecom Italia dataset [6]. Finally, we conclude the article.

BACKGROUND

5G CORE NETWORK ARCHITECTURE

The 5G CN architecture has been devised with the aims of maintaining compliance with the concept of 3GPP communication reference interfaces (point-to-point), and evolving in the near future to a service-oriented architecture, where the 5G core control network functions communicate via a common bus, through the concept of producers and consumers. Indeed, the 3GPP SA2 group has defined two architectures, one based on the reference interfaces and the other following a service oriented architecture. In this article, we focus on the second one as the new components are expected to be run on top of a virtualized environment. For the first architecture, readers may refer to [2]. Regarding the 5G core functions, some are similar to 4G Evolved Packet Core (EPC) ones, and some are completely newly designed. Notably, the access control and session management functions are combined in the mobility management element (MME) of EPC, but separated in the 5G core to better support fixed access, and ensure scalability and flexibility. The most relevant network functions, as defined in 5G Core, are:

- Core AMF handles access control and mobility. If fixed access is involved, the mobility management will not be required in the AMF. It is in charge of handling attach requests of UEs from several RANs; as the entry point to the CN control plane, it creates bottleneck issues when undersized.

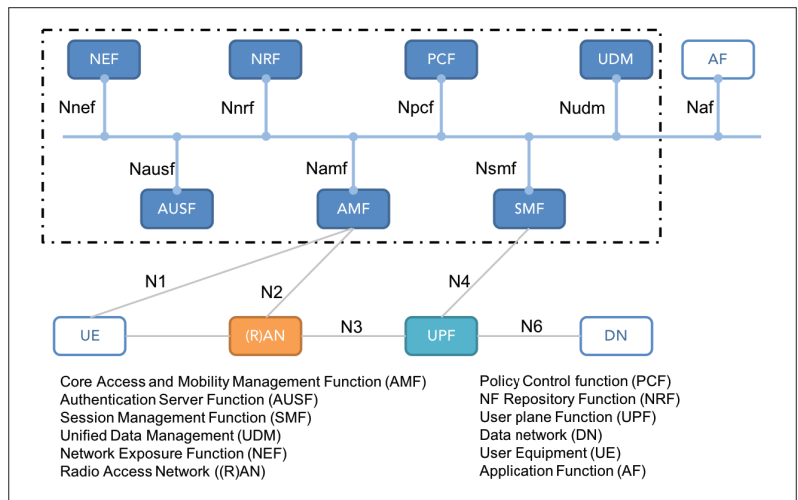


FIGURE 1. Service oriented 5G core network.

- Session management function (SMF) sets up and handles sessions according to the network's policy.
- Policy control function (PCF) corresponds to the policy and charging rules function (PCRF) in LTE. It is worth noting that this function will also integrate a policy framework for network slicing.
- User plane function (UPF) can be deployed, based on the service type, in several configurations and locations to be able to handle and forward users' data traffic.
- Unified data management (UDM) is similar to the home subscriber server (HSS) in EPS. However, it is envisioned that it integrates subscriber information for both fixed and mobile access in the new generation (NG) core.
- NF repository function (NRF) is a new function. It provides registration and discovery functionality, allowing the CN functions to discover each other and communicate via application programming interfaces (APIs).

In the following we briefly describe the AMF and SMF functions as this work focuses on the AMF scalability. For more details on the other functions, the reader may refer to [2]. The AMF ensures UE-based authentication, authorization, mobility management, and so on. UEs, even using multiple access technologies, are connected to a single AMF, since the latter is agnostic to the access technology. The SMF is responsible for session management and IP address allocation to UEs. In addition, it selects and controls the UPF for data transfer. When a UE has multiple sessions, several SMFs may be allocated to each session, aiming at managing the latter individually and possibly providing different functionalities per session. Based on operator policy or AF provided policy on packet flow, the PCF defines policies about mobility and session management, which will be enforced by the AMF and SMF.

Figure 1 illustrates the service oriented 5G CN architecture. Only the UPF is deployed on the legacy transport network infrastructure; the other CN functions handling the control plane are expected to run on the more centralized telco cloud platform. While the communications with the UE, RAN, and UPF use the reference interfaces, the

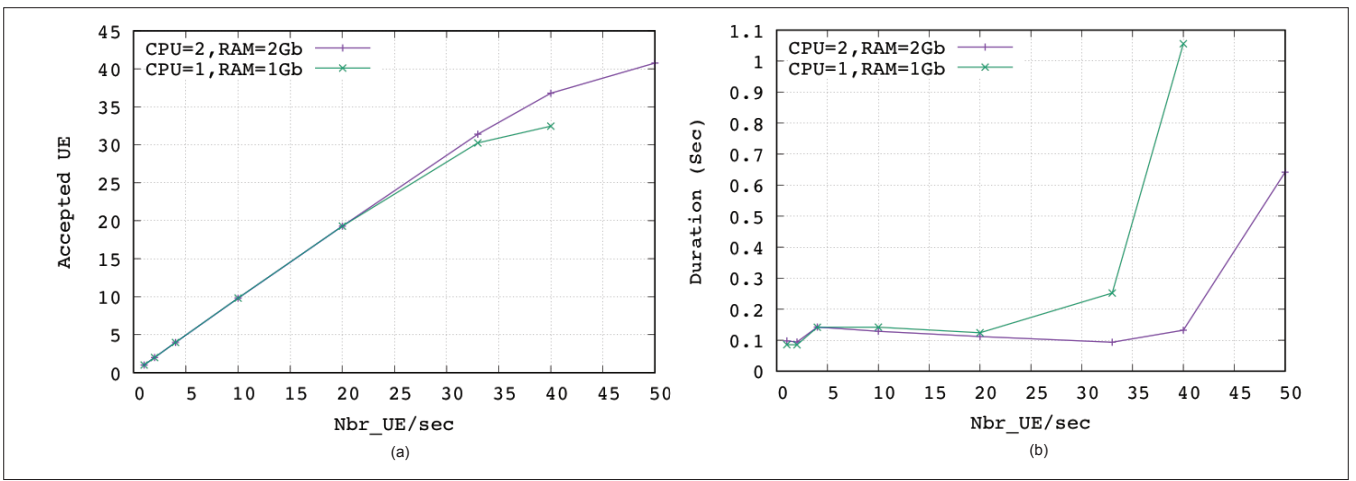


FIGURE 2. Benchmarking result of the AMF: a) accepted UEs vs. the number of UE attach requests; b) attach duration vs. the number of UE attach requests.

core control plane functions are connected via a common communication bus. As stated earlier, the CN functions use the concept of producer/consumer, where a CN function (CNF) may register, using the NRF, for specific events provided by another CNF via an API. Thus, either 1:1 or 1:N communication is possible. For instance, the AMF service exposes information to the other CNFs regarding events and statistics related to mobility, or a PCF provides all the operations related to policy rules to other CNFs. The service-oriented 5G core architecture could easily be deployed in a virtualized environment (e.g., VM or container), wherein libraries of functions may be requested from a VNF catalog and composed into a virtual CN on demand.

RELATED WORK

As illustrated in Fig. 1, the AMF is the entry point of the 5G CN, which makes it the system bottleneck. Thus, special focus should be given to the dimensioning of the VNF(s) hosting the AMF to ensure system scalability. Scalability in this case is more a software deployment issue than a hardware upgrade as most of the CNFs are running as VNFs on top of a virtualized infrastructure.

One of the most common scaling solutions consists of using static thresholds. In [7–9] the authors have proposed scalability based on thresholds. In the same spirit, we also worked on threshold-based scalability in a previous attempt to manage AMF instances in [5]. The main idea is to deploy a new instance (scale-out) when the CPU usage (or other resources) exceeds a certain fixed “scale-out” threshold. Reciprocally, an instance is shut off (scale-in) when a considered threshold exceeds a fixed value. Static-threshold-based solutions are already deployed and widely used in commercial and open source cloud solutions, like [10, 11]. However, static-threshold-based approaches are not suitable for mobile networks, and have several limitations due to the fact that they react to events (i.e., reactive solution). Indeed, deploying a new instance of a VNF when reacting to events may take from one minute to a dozen minutes depending on the data center architecture and the complexity of scaling the VNF instance, leading to delay in the

scaling-out process, hence increasing the attach duration of UE and rejection of some requests.

An alternative to static-threshold-based solutions is the usage of adaptive techniques, as proposed in [12]. In this work, the authors proposed a mechanism that combines Q-learning with Gaussian-process-based system models, which allows adaptation to dynamic environments and improves the scaling policy before taking any action. The authors assume that their proposal reacts better than static threshold rules. However, it remains a reactive solution, and hence inherits the same weaknesses. Moreover, using reinforcement learning, it may take a considerable time before starting to have the correct decisions, which is not acceptable in systems like 5G due to the time constraint of certain types of applications. Note also the side effects associated with this type of approach, consisting of the influence of the external environment on performance and therefore the risk of quality degradation.

Proactive solutions for scalability have also been proposed in the literature. The authors of [13] suggested to proactively schedule resources for VNFs based on CPU usage forecasting from a historical dataset. Although this technique is interesting for its proactivity, it does not offer adaptation in case of traffic pattern evolution and is not compliant with mobile network requirements. Indeed, the proposed solution considers only system-level information, like consumed CPU or memory, ignoring service-level information. VNF scaling in mobile networks should take into account strict rules as defined in standards; therefore, deploying only a new instance by a cloud orchestrator is not sufficient for scaling an end-to-end service. Consequently, mobile networks require service-level information, like the number of connected users or traffic load when scaling-out or -in. Moreover, the parameters of the “offline” prediction equation proposed by the authors are determined statically and are completely independent of the considered dataset. This implies stability in the weights that are given for daily prediction, and thus a risk of divergence in the prediction with the slightest change of pattern. Concerning the “online” approach, the authors propose to use an average (exponen-

tial average) estimator. However, following [14], when a pattern can be determined, this is clearly not the best solution.

ML-BASED SCALABLE MECHANISM FOR 5G CORE NETWORK

SCALABILITY ISSUE OF THE AMF

To illustrate the need for an efficient scalability mechanism for the AMF, we conducted a set of experimentations using OpenAirInterface (OAI: <http://www.openairinterface.org/>, accessed March 2018). OAI is a set of open source software tools that allows rapid prototyping of 5G mobile network, including the RAN and CN functions. We used the emulated version of OAI, known as OAI simulator (OAISIM), which enables the emulation of a set of connected UEs and one eNodeB. The AMF was run on top of a VM, with two configurations: 1 CPU and 1 GB of memory (noted as the first configuration); 2 CPUs and 2 GB of memory (noted as the second configuration). It is important to note that we refer to AMF, even if we used the MME function of the OAI EPC. Indeed, our focus is mainly to compute the number of accepted UEs and the time needed for a UE to attach to the network when varying the number of attach requests per second. Both metrics refer to the performance of an authentication procedure, which is one of the common functions between the MME and AMF.

Figure 2a represents the accepted UEs according to the number of UE attach requests per second. Clearly, we observe that from a certain threshold, both configurations start rejecting UE attach requests. Obviously, the stronger the configuration of the VM, the higher the threshold. Using the first configuration, the threshold is 20 UEs/s, while for the other configuration is around 35 UEs/s. We remark the same trend regarding the time needed by a UE to successfully attach to the network, which is presented in Fig. 2b. In this case, we observe that from the same threshold, 20 UEs/s for the first configuration, and 40 UEs/s for the second one, the attach duration is exponentially increasing, which violates the requirements specified by the 3GPP standard [15]. Obviously, there is a need to have a mechanism to solve this issue, which may consist of either adding more resources (CPU and memory) to the VM, that is, scaling-up, or adding a new VM to the AMF function (with an internal load balancer), that is, scaling-out.

TRAFFIC FORECASTING: DEEP-LEARNING-BASED APPROACH

As explained above, there is a need for a mechanism to allow dynamic dimensioning of the AMF resources to avoid degrading the system QoS. The envisioned mechanism should be proactive to avoid the weakness of reactive solutions, hence requiring a traffic prediction solution. For this aim (i.e., traffic prediction), two neural-network-based techniques were envisioned and tested, aiming to predict the upcoming traffic load. The first one is based on a deep learning neural network (DNN), while the second one relies on a recurrent neural network (RNN), more specifically on a long short term memory (LSTM) cell. In a DNN, a particular neuron can be activated by the contribution of

The main procedures for traffic prediction using learning techniques are: getting the dataset, formatting the data, designing the neural network, training the neural network and then starting predicting (or testing your network).

several inputs at once, and therefore there is no real recognition of a traffic pattern. On the contrary, an RNN, and particularly an LSTM network, receives the inputs successively, and through a memory and omission mechanism, it manages to memorize traffic patterns, as has been demonstrated in several works [16].

The main procedures for traffic prediction using learning techniques are: getting the dataset, formatting the data, designing the neural network, training the neural network, and then starting prediction (or testing your network). Each step is important in order to get the best accuracy of prediction compared to real values.

Dataset: The main key of prediction using deep learning is the dataset. Nowadays, operators are able to collect those data easily as they are the owner of their infrastructure and have a clear idea of the evolution of the load over their network. In this work, we have used the dataset collected during the Big Data Challenge organized by Telecom Italia [17]. The dataset is open source [18] and rich in information. It can be used to predict the evolution of the load in a CN.

The considered dataset contains several kinds of data, like sms in/out, call in/out, and the call detail records (CDRs) with an interval of 10 minutes for two months, including information on both control plane and data plane traffic. We consider in this work that at least 10 percent of the traffic is control plane traffic, and the remaining 90 percent is data plane traffic. Therefore, for the remainder of this article, we only use the control plane traffic (10 percent of the load in each period of 10 minutes) for prediction.

Data Formatting: Once the dataset is collected and 10 percent of each period is extracted, the data should be formatted in order to be injected in the neural network for the training phase. This part requires close attention as the data format can affect the accuracy of prediction of the neural network. Intelligently formatting the data will allow the neural network to learn more in the training phase, and thus predict with high accuracy. However, adding lots of key performance indicators (KPIs) may lead to decreasing the prediction accuracy further, particularly for small datasets. Therefore, a balance should be found between the dataset size and the KPIs injected in the data formatting. In our case, we decided to format the dataset by classifying the load into 10 different classes, which can be interpreted as: class 1 is a load interval where (x) AMF instances are needed; Class 2 (x + y) AMF are needed, and so on. Note that adding a much larger number of output classes would require a much larger dataset. Indeed, the addition of classes implies the increase of the number of parameters of the neural network, which has the effect of requiring more data to have a significant precision.

Neural-Network-Based Prediction: The first technique used for predicting the class of load of the next period is DNN. DNN is a neural network

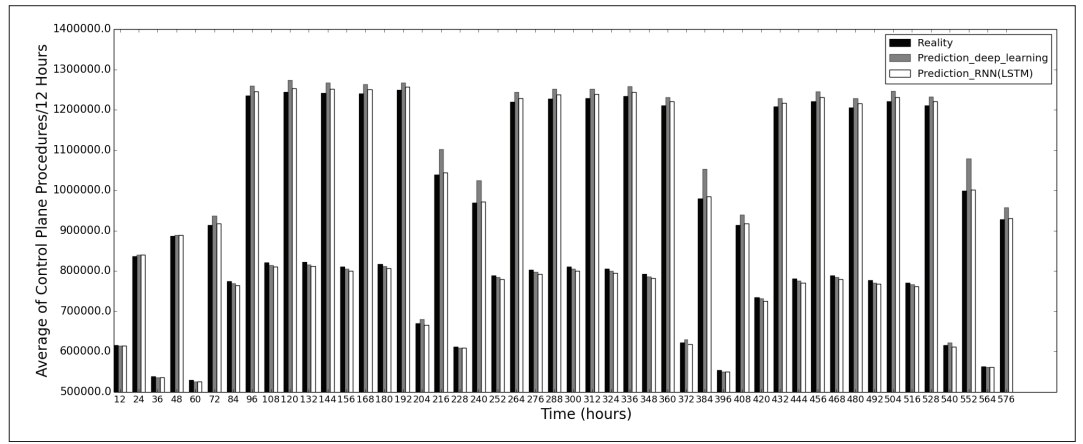


FIGURE 3. Reality vs. NN prediction.

There is a need for a mechanism to allow dynamic dimensioning of the AMF resources to avoid degrading the system QoS. The envisioned mechanism should be proactive to avoid the weakness of reactive solutions, hence requiring a traffic prediction solution.

composed of three main layers: the input layer, the middle layer(s) (hidden layer(s)), and finally, the output layer. This type of neural network is also known as the feed forward neural network (FFNN). The data goes in the neural network from the input layer through the middle layer(s), and finally go out through the output layer.

The second technique tested for predicting the average load of the upcoming period of time is the RNN, which is another type of neural network. Unlike the DNN, the RNN has loops. Those loops allow injecting previous events while predicting future events. It is a sort of memory block allowing decisions to be made depending on the previous event(s). Also, in an RNN multiple kinds of cells are used today. In this article we choose to rely on the LSTM cell. The advantage of the LSTM cell is to allow data patterns to be stored without degradation over time.

In order to compare the accuracy of the DNN and RNN (LSTM), one DNN and one RNN (LSTM) were designed using the Tensor Flow library [19]. We divided and formatted the data into classes of average load. Then we split the original dataset into two separate datasets. The first one contains 60 percent of the data for training the two neural networks, while the remaining 40 percent are left to test the robustness of the two networks. Both networks are trained with the 60 percent dataset and then asked to predict the remaining 40 percent. Figure 3 illustrates the prediction of the DNN and the RNN. From Fig. 3 we notice that the RNN (LSTM) performs better than the DNN. The RNN (LSTM) achieves an accuracy of ≈ 90 percent, while the DDN obtains ≈ 80 percent of accuracy; hence, the RNN performs more than 10 percent better than the DNN when predicting the class of average load of the upcoming period of time. The accuracy of the prediction is mainly due to the ability of LSTM networks to recognize a traffic pattern, unlike DNNs. The accuracy of the DNN can be improved, but this would require more neurons and more data, something that is not always available.

Based on these results, we consider using the RNN (LSTM) for the AMF scalability mechanism.

PREDICTION-BASED SCALABILITY

Having described the prediction model, now we present the prediction-based scalability algorithm. Based on the ETSI network functions virtualization (NFV) reference architecture, the scalability of a VNF is a decision made by the NFV orchestrator (NFVO), using the information provided by both the virtual infrastructure manager (VIM) and the VNF manager (VNFM). The former provides information on the VM execution environment (e.g., CPU usage, memory usage), while the latter provides information obtained from the element manager (EM) hosted in the VNF. The information provided by the EM represents the inside VNF service performance, like the number of attach requests or the size of the queue of the AMF. In [4], we proposed that the EM of the AMF reports this information to the VNFM, hence the NFVO. The latter makes decisions to scale-in or -out using predefined thresholds. However, due to the time needed to scale out (run a new VM), the proposed solution suffers from delay in reacting to the increase of UE attach, which leads to an increase in the time of users' attachments with some loss of requests.

To remedy this problem, in this work, we use the RNN to predict the needed number of AMFs that efficiently manages the traffic load, aiming at keeping the attach duration under an acceptable value (using the previous results). The NFVO uses the RNN to predict the traffic (i.e., class of traffic), and hence derives the optimal number of AMFs to proactively deploy before the change of traffic load, anticipating any network load changes and removing the time needed to start a new VM instance in the scale-out process. Besides, in this work there is no need to monitor the VNF performance via the VNFM as the NFVO relies only on the trained RNN, which permits reduction of the communication load introduced by the exchanged messages between the EM and VNFM.

PERFORMANCE EVALUATION

To evaluate the performance of the proposed solution and compare it with the threshold-based mechanism, we developed a discrete event simulator using Simpy [20]. The choice of Simpy

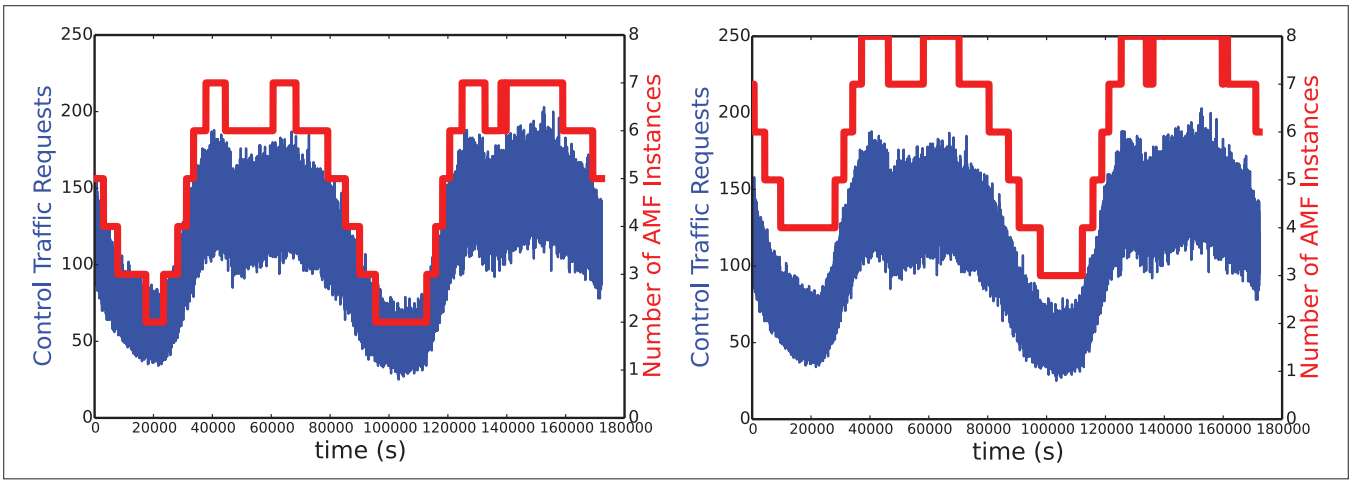


FIGURE 4. Evolution of AMF instances number with control traffic requests: a) threshold-based solution; b) RNN-based solution.

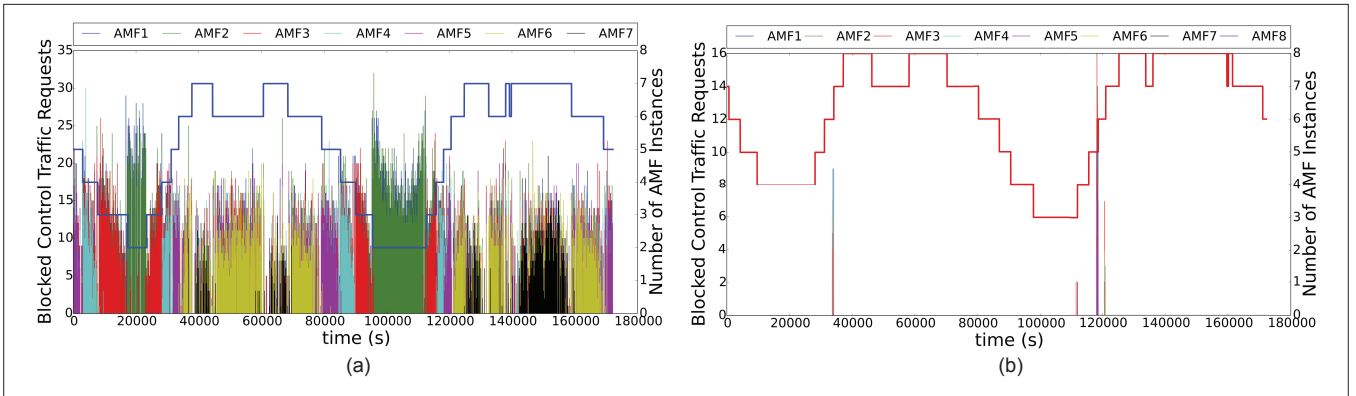


FIGURE 5. Number of blocked control traffic requests compared to the number of AMF instances: a) threshold-based solution; b) RNN-based solution.

was motivated by its efficiency when it comes to developing a simulator from scratch as well as its integration with other tools we have used, such as tensorflow. The simulator consists of the following:

- A UE traffic generator allows initiation of UE attach requests according to the dataset.
- An eNodeB dispatches the arrival of control requests upon the deployed instances of AMFs using a load balancing mechanism if more than one AMF instance is used. It is based on the algorithm introduced in [5].
- Each VNF (hosting an AMF component) can process up to 20 UE attach requests per second.
- An NFVO manages the scalability of the AMF.

We launched several simulations to compare the performance of the proposed RNN-based scalability mechanism with the threshold-based solution. The former predicts the average load for the next period of time (10 minutes) and deploys the exact number of AMF instances needed to handle the upcoming load, while the latter launches a new instance (scale-out) when the overall arrival average load exceeds 80 percent of the overall system capacity (depending on the number of active AMF instances at the moment), and removes an instance (scale-in) when the overall average load decreases below 10 percent of the overall CN capacity. Both mechanisms were

tested over 48 hours of arrival data and with a deployment latency (i.e., time needed for the AMF instance to be operational) of 60 seconds. We chose this value according to [21], wherein the authors have studied the deployment latency for different VNF deployment scenarios on top of a private cloud infrastructure.

Figure 4 illustrates the attach request rates vs. the number of AMF instances deployed by both mechanisms; that is, threshold-based and RNN-based. From Fig. 4 we notice that both mechanisms are able to follow the arrival flow by deploying more AMF instances when the arrival load increases and reducing the number of AMF instances when the traffic load decreases. We remark that the RNN-based mechanism follows the arrival load perfectly. Furthermore, the RNN-based solution always anticipates the traffic load, and the number of AMF instances always follows the traffic trends. In contrast, the threshold-based mechanism takes more time to react. Indeed, we observe that in many cases, the two curves (traffic load and number of AMFs in Fig. 4a) cross each other, meaning that the threshold-based solution is taking time to react, leading to degraded network QoS.

This observation is confirmed in Fig. 5, which illustrates the number of blocked attach requests in both solutions. Indeed, we remark that several requests are rejected in the case of the threshold-based solution compared to the RNN-based

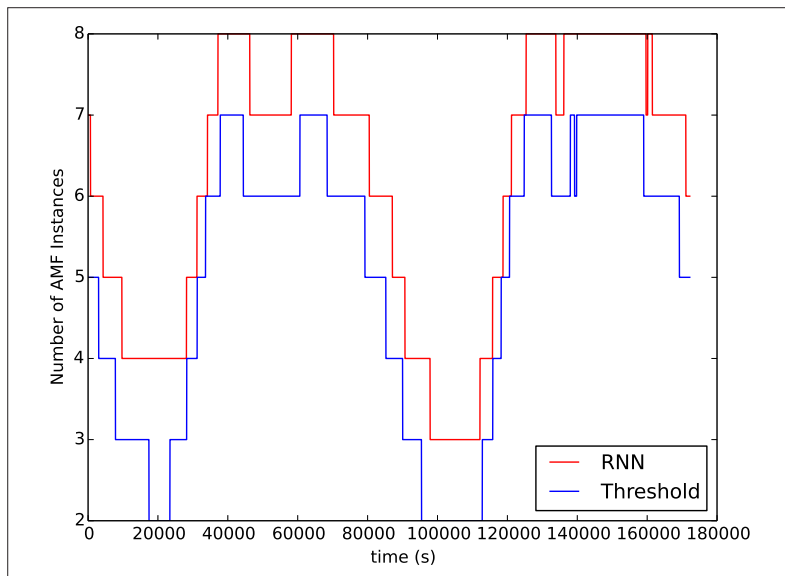


FIGURE 6. Number of deployed AMF instances (threshold scaling vs. RNN scaling).

solution. The accurate prediction of the RNN-based solution allows reacting in advance of the traffic changes by proactively instantiating the necessary VNF resources, avoiding overloading the AMF while waiting for the additional resource to be instantiated. Moreover, this result confirms that the RNN-based solution accurately estimates the needed number of AMF instances to manage a specific traffic load. For example, when the load is maximal, the RNN-based mechanism deploys up to 8 instances, while the threshold-based solution deploys up to 7. Reciprocally, for the case when the arrival load is minimum, the threshold-based mechanism deploys at least 2 instances, while the RNN-based mechanism instantiates at least 3 AMF instances. Thus, the RNN calculates more accurately the number of instances needed to manage the load, allowing the system to be stabilized and significantly decreasing the number of rejected requests.

Figure 6 illustrates the number of deployed AMF instances using the threshold-based mechanism vs. the RNN-based mechanism. This figure validates the trend seen in the preceding figures, where we clearly observe that the RNN-based mechanism is always in advance compared to the threshold-based solution while reacting to traffic change. Thanks to traffic prediction, the deployment delays are removed, particularly for the scale-out process, which is not the case for the threshold-based mechanism, which suffers from latency while deploying a new instance. Indeed, as the RNN predicts the arrival load of the upcoming period, it deploys the AMF instances needed before the upcoming load in such a way that the AMF instances are up and ready with the increase of the load.

To sum up, the RNN(LSTM) technique performs better than the threshold technique, first because it acts in a proactive way unlike the threshold technique. Second, it allows accurate deployment of the exact needed number of AMF instances to absorb the upcoming load. Finally, it allows removal of the deployment latency by scaling out in advance the needed AMF instances. In this way, the mobile core is ready to absorb the

upcoming load increase without suffering from overhead and extra latency for the control procedures.

CONCLUSION AND FUTURE WORK

In this article, we propose a novel solution, based on machine learning, to scale-out and -in the AMF resources in a virtualized environment. Using a neural network that has been trained on a mobile network traffic dataset to forecast the user attach request rate allows prediction of the exact number of AMF instances needed to handle the upcoming user traffic. By being proactive, the proposed solution allows the deployment latency to be avoided when scaling out resources. The latter may degrade the network performance, since the AMF is overloaded while waiting for the additional resource to be instantiated. Simulation results confirm the efficiency of the RNN-based solution compared to a threshold-based solution. Future work aims to use the trained RNN to estimate the number of user plane 5G CN functions needed to handle the data plane traffic in a virtualized environment.

ACKNOWLEDGMENT

EURECOM's contribution has been partially funded by the European Framework Program under H2020 grant agreement No. 723172, 5G!Pagoda project.

REFERENCES

- [1] ETSI, "Network Functions Virtualisation (NFV); Architectural Framework"; <https://bit.ly/2GGPZpN>, accessed Mar. 2018.
- [2] 3GPP, "System Architecture for the 5G System"; <https://bit.ly/2GM9zgo>.
- [3] 3GPP, "Telecommunication Management; Study on Management and Orchestration of Network Slicing for Next Generation Network"; <https://bit.ly/2qeZd1s>.
- [4] I. Alawe et al., "On Evaluating Different Trends for Virtualized and SDN-Ready Mobile Network," *Proc. 2017 IEEE 6th Int'l. Conf. Cloud Networking*, 2017, pp. 1–6.
- [5] I. Alawe et al., "On the Scalability of 5G Core Network: The AMF Case," *Proc. IEEE Consumer Commun. and Networking Conf.*, 2018.
- [6] Telecom Italia, Telecom Italia Big Data Challenge; <http://www.telecomitalia.com/tit/en/bigdatachallenge/contest.html>, accessed Mar. 2018.
- [7] S. Dutta, T. Taleb, and A. Ksentini, "QoE-Aware Elasticity Support in Cloud-Native 5G Systems," *Proc. IEEE ICC*, 2016, pp. 1–6.
- [8] G. A. Carella et al., "An Extensible Autoscaling Engine (AE) for Software-Based Network Functions," *IEEE NFV-SDN*, 2016, pp. 219–25.
- [9] A. B. Alvi, T. Masood, and U. Mehboob, "Load Based Automatic Scaling in Virtual IP Based Multimedia Subsystem," *Proc. IEEE CCNC*, 2017, pp. 665–70.
- [10] Amazon, "Web Services Auto Scaling"; <https://aws.amazon.com/autoscaling/>.
- [11] A. CloudStack, "Open Source Cloud Computing"; <http://cloudstack.apache.org/>.
- [12] C. H. T. Arteaga, F. Rissio, and O. M. C. Rendon, "An Adaptive Scaling Mechanism for Managing Performance Variations in Network Functions Virtualization: A Case Study in an NFV-Based EPC," *Proc. 2017 IEEE 13th Int'l. Conf. Network and Service Management*, 2017, pp. 1–7.
- [13] A. Bilal et al., "Dynamic Cloud Resource Scheduling in Virtualized 5G Mobile Systems," *Proc. 2016 IEEE GLOBECOM*, 2016, pp. 1–6.
- [14] R. Fu, Z. Zhang, and L. Li, "Using LSTM and GRU Neural Network Methods for Traffic Flow Prediction," *Proc. IEEE Youth Academic Annual Conf. Chinese Association of Automation*, 2016, pp. 324–28.
- [15] 3GPP, "Non-Access-Stratum (NAS) Protocol for Evolved Packet System (EPS); Stage 3"; <https://bit.ly/2H0Tred>, accessed Mar. 2018.
- [16] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, 1997, pp. 1735–80; <https://doi.org/10.1162/neco.1997.9.8.1735>.

- [17] 3GPP, "Non-Access-Stratum (NAS) Protocol for Evolved Packet System (EPS); Stage 3"; <https://bit.ly/2H0Tred>, accessed Mar. 2018.
- [18] T. Italia, "A Multi-Source Dataset of Urban Life in the City of Milan and the Province of Trentino Dataverse"; <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/EGZHFV>.
- [19] "TensorFlow, An Open-Source Machine Learning Framework for Everyone"; <https://www.tensorflow.org/>, accessed Mar. 2018.
- [20] 6Team-SimPy, Discrete Event Simulation for Python; <https://simpy.readthedocs.io/en/latest/>, accessed Mar. 2018.
- [21] P. A. Frangoudis, L. Yala, and A. Ksentini, "CDN-as-a-Service Provision over a Telecom Operators Cloud," *IEEE Trans. Network and Service Management*, vol. 14, no. 3, 2017, pp. 702–16.

BIOGRAPHIES

IMAD ALAWE is an engineering graduate of the Rennes School of Engineering and a Ph.D. student at TDF and IRT b<>com. He started at IRT b<>com in 2015. He first worked on SDN for optical networks in collaboration with Ekinops; then he joined TDF to work toward a Ph.D. He works on the integration of SDN in the EPC, with NFV and architecture evolutions toward 5G architecture.

ADLEN KSENTINI is an IEEE Communications Society Distinguished Lecturer. He obtained his Ph.D. degree in computer science from the University of Cergy-Pontoise in 2005. From 2006 to 2016, he was an assistant professor at the University of Rennes 1. In March 2016, he joined the Communication Systems Department of EURECOM as an assistant professor. He has been working on network slicing and 5G in the context of two European projects on 5G, H2020 projects 5G!Pagoda and 5GTransformer.

YASSINE HADJADJ-AOUL is working as an associate professor at the University of Rennes, France, where he is a member of the IRISA Lab and the INRIA Dionysos project team. He received a Ph.D. degree in computer science from the University of Versailles, France, in 2007. He was a postdoctoral fellow at the University of Lille and University College Dublin under the EUPF6 Marie Curie Action. His main research interests concern the fields of congestion control and QoS provisioning.

PHILIPPE BERTIN is a senior research engineer at Orange Labs. He is managing research projects on networks architecture with b<>com Institute for Research and Technology. His research interests include the design of distributed and dynamic control and data planes for flexible and convergent 5G networks, leveraging on SDN and virtualization. He graduated from Paris 6 University (M.Sc., 1993) and Rennes University (Ph.D. in computer science, 2010).