



Autonomous Network Slicing Prototype Using Machine-Learning-Based Forecasting for Radio Resources

Nazih Salhab, Rami Langar, Rana Rahim, Sylvain Cherrier, Abdelkader Outtagarts

► To cite this version:

Nazih Salhab, Rami Langar, Rana Rahim, Sylvain Cherrier, Abdelkader Outtagarts. Autonomous Network Slicing Prototype Using Machine-Learning-Based Forecasting for Radio Resources. IEEE Communications Magazine, 2021, 59 (6), pp.73-79. 10.1109/MCOM.001.2000922 . hal-04049769

HAL Id: hal-04049769

<https://univ-eiffel.hal.science/hal-04049769>

Submitted on 28 Apr 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Autonomous Network Slicing Prototype Using Machine Learning Based Forecasting for Radio Resources

Nazih Salhab, Rami Langar, Rana Rahim, Sylvain Cherrier and Abdelkader Outtagarts

Abstract—With the emergence of virtualization and software automation for mobile networks, network slicing is enabling operators to dynamically provision network resources tuned to suit heterogeneous service requirements. This article investigates the architectures of the Fifth Generation (5G) of mobile networks experimental prototypes with a focus on network slicing. We present some existing 5G prototypes and identify their gaps. We, then, propose an architecture and a design of a 5G micro-service based prototype. Such prototype has an ability to auto-configure radio resources for network slices using machine learning (ML) powered decisions based on real-time acquired performance metrics. Finally, we discuss some use cases on top of this prototype and their related results before concluding.

Index Terms—Network Slicing prototype, micro-service architecture, OpenAirInterface (OAI), TICK Stack

I. INTRODUCTION

ON top of conventional communication scenarios among people, the Fifth Generation (5G) of mobile communications and Beyond (B5G) will empower machine-to-machine communication that drive different applications for multiple industries. Vertical applications with heterogeneous requirements in terms of latency, reliability, mobility support, energy efficiency, throughput, and connection/traffic densities will flourish with 5G/B5G. These applications include smart grid, collaborative robots, cooperative vehicles, smart farms, and so on. To serve such a plethora of applications, mobile network operators need to modernize their architectures to accommodate forecasted exponential increase in terms of number of connected devices and explosive size of exchanged data. Networks have to become agile by offering tailor made infrastructures to suit specific use-cases and their requirements. 5G Standards Developing Organizations envision supporting at least three major categories of services: enhanced Mobile Broadband (eMBB), massive Machine-Type Communications (mMTCs), and ultra-Reliable Low-Latency Communications (uRLLCs). These are considered as network slices consisting of end-to-end logical infrastructures deployed over a shared physical infrastructure inline with sought performance demands. Network slicing and dynamic provisioning are especially interesting in case of uRLLC, where virtualized network functions have to be deployed at the edge of the network to minimize the latency. Conversely, traditional system architectures rely on monolithic systems. In the monoliths, the presentation layer (user interface) and the data access layer (database) are tightly coupled in a single program making the application self-contained and independent from other

computing applications, but also, highly dependent on hosting platform. On the other hand, monoliths have some downsides, especially, when it comes to the lack of flexibility in terms of maintenance complexity and evolution speed. With data usage growth, monoliths become too large to operate and maintain. In addition, in order to upgrade a functionality in a particular monolith, system administrators have to take the whole system down to complete such an operation. Moreover, a risk of regression during integration that causes unplanned delays should not be neglected. Therefore, traditional monoliths cannot provide required agility for 5G/B5G. Advents of Network Function Virtualization (NFV) and Software-Defined Networking (SDN) allow overcoming this limitation [1].

On another front, the cloud positions itself as an enabler for 5G/B5G and is well suited, in particular, for the Radio Access Network (RAN) as Cloud-RAN (C-RAN) with inherent elasticity. Instead of relying on dedicated servers and proprietary computing devices, the cloud provides virtual, but dedicated resources relying on shared capabilities (Compute, Networking and Storage). Cloud resources provide required flexibility, featured through auto-scalability, and reduced deployment time, while ensuring up-to-date software releases [2].

In this context, we propose, in this paper, a micro-service based experimental prototype with Machine Learning (ML) forecasting capability for autonomous slicing management. Based on an exhaustive evaluation of multiple ML techniques [3], we implement the best performer prediction method and leverage our prototype to analyze its performance in a closed-loop setting. Also, compared to our former prototype [3], we consider, in this paper, a micro-service based deployment of the core network rather than a monolithic one. Moreover, we have developed special scripts for the data processing engine and interfaced the latter with a northbound application for an autonomous configuration management.

The rest of this article is organized as follows. We first investigate the state of the art, the major European projects and the standardized reference architecture for 5G. Then, we present our approach in designing our 5G experimental prototype with machine-learning capability using open-source software. Finally, we present two use cases of the designed prototype using our Internet of Thing (IoT) programming cloud platform, called BeC3 (Behavior Crowd Centric Composition) [4] and our ML based autonomous slicing prototype with forecasting capability and discuss related experimental results.

II. EXISTING 5G PROTOTYPES

Authors in [1] surveyed 5G Network slicing using SDN and NFV. They discussed the related architectures and future challenges. They clearly advocated containerization technologies as an enabler for 5G. Moreover, author in [2] discussed the pros and cons of cloud migration and the points in between, suggesting seven R's (Replace, Reuse, Refactor, Re-platform, Re-host, Retain and Retire) strategy, while emphasizing on cloud architecture advantages. They addressed agnostic applications migration issues. Therefore, we proposed a microservice based cloud native architecture for our prototype using Docker [5].

Authors in [6] introduced 5G-EmPOWER platform that enables complex policies deployment and management over SDN. They evaluated the performance of their platform using Long Term Evolution (LTE) network elements provided by Software Radio Systems (srsLTE).

Authors in [7] presented a prototype of a virtualized 5G infrastructure supporting network slicing. They investigated how OpenAirInterface™ (OAI) [8] can be used on top of Mobile Central Office Re-architected as a Datacenter (M-CORD). Moreover, they demonstrated a virtualized 5G RAN and core deployment, using monolithic architecture. Recall that OAI is an open-source software implementation of 4G/5G, spanning the full protocol stack of both radio and core networks. On the other hand, M-CORD is an open source reference solution for carriers deploying 5G mobile wireless networks.

The O-RAN [9] Alliance initiated an operator driven mission to open the RAN interfaces in an aim to allow multi-vendor equipment to work seamlessly across well-defined interfaces. O-RAN aims to transform the RAN Industry towards open, intelligent, virtualized and fully interoperable one.

MOSAIC-5G [10] encompasses multiple projects in an aim to transform the RAN and the core into an agile service delivery platform to rapidly explore new ideas as well as emergent application and changing business needs. In particular, we have reused the FlexRAN project [10], [11] to implement the RAN SDN controller and automate its control through our northbound web app for an autonomous configuration management.

Authors in [12] reported their experience building a 5G slicing prototype that enables multi-tenancy through virtualization. They highlighted the lack of monitoring capabilities, cross-domain and intelligent orchestration for such prototypes and flagged these axes as open research questions. We note that these works did not include a holistic performance management stack or illustrate its benefits through common use cases. Therefore, we identify four objectives that a 5G prototype needs to address: i) minimizing the total cost of ownership by relying on open source software, ii) automating reconfiguration through data-driven ML decisions, iii) overcoming IoT complexity through abstraction, and iv) minimizing human management by employing self-healing, resilience and auto-scaling. Our proposed 5G prototype will indeed provide these features, as described later.

III. EUROPEAN PROJECTS AND STANDARDS DEVELOPING ORGANIZATIONS

The 5G Infrastructure Public Private Partnership (5G PPP) includes several projects aiming to formalize the reference architecture of 5G. Some of these projects are: i) 5G-NORMA (Novel Radio Multiservice adaptive network Architecture), ii) 5G-xHAUL (5G-Cross Haul), iii) COHERENT (Coordinated control and spectrum management for 5G heterogeneous radio access networks), iv) Euro-5G, v) Fantastic 5G (Flexible Air Interface for Scalable service delivery within the wireless communication networks of the 5G), vi) Flex5Gware (Flexible hardware/software platforms for 5G network elements and devices), vii) METIS-II (Mobile and Wireless communications Enablers) and viii) 5G-ESSENCE (Edge Cloud computing and Small Cell as a Service). On the other hand, several standards developing organizations such as 3rd Generation Partnership Project 3GPP, ETSI NFV evolution and ecosystem, Next Generation Mobile Networks (NGMN), TeleManagement Forum (TM Forum) and Open Networking Foundation (ONF) have envisioned a service based architecture for 5G.

IV. ARCHITECTURE OF OUR PROPOSED PROTOTYPE

We propose using microservice based architecture that employs a collection of autonomous self-contained services, each implementing a single business capability. It is an evolution of a service-oriented architecture by further decomposing the services into smaller, independent microservices with loose coupling relation. A well-known example of this architecture is a cluster of microservices managed by an orchestrator. This architecture is convenient for large applications with numerous sub-domains that require high release velocity and high scalability. However, there are some downsides for this architecture, that are: i) accrued latency due to distribution of microservices over several machines and ii) potential internal network congestion as more inter-service communication is involved. Note that such latency can be minimized by setting up a local repository to store the different network function images. Moreover, the high latency drawback appears when a virtualized network function is deployed from scratch. In contrast, when a scaling operation takes place, the change in terms of microservices replica happens seamlessly, without affecting the available microservices. In the following, we detail our proposed architecture for 5G experimental prototype as depicted in Fig. 1.

A. Infrastructure Layer (5G system)

To be in line with the reference 5G architecture, we deploy OpenAirInterface for Core-Network (OAI-CN) [8] as Docker containers in a Docker Swarm Cluster consisting of a manager and two workers. Note that we propose using Docker Swarm instead of Kubernetes [14] due to its native support by Docker, being the default orchestrator offered by Docker and thus, it is much simpler when it comes to implementation in Lab testing environments. Moreover, using containers allows us to manage the infrastructure in an agnostic way allowing us to seamlessly upgrade it to different 5G releases whenever they are ready. Note also that we considered a number of replicas

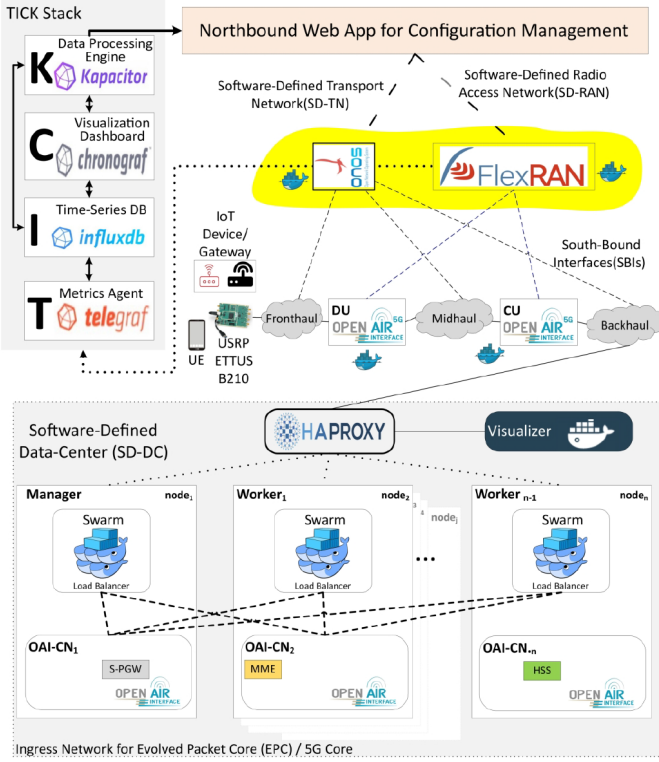


Fig. 1. Experimental Prototype Architecture

that is strictly greater than zero.

Accordingly, using the Docker Swarm container orchestration tool allowed us to seamlessly scale the number of replicas of each microservice without impacting the service to which they belong. Therefore, we have included in our auto-scaler a stopper-condition to enforce this, in order to avoid any downtime. However, in case of a failure of the hosting docker-machine, and supposing that the whole microservices are affected, it takes multiple seconds to redeploy the corresponding microservices. The downtime depends on multiple factors including loading it from the cached Docker images or having to download it from the online or local repositories. We also implement the next generation fronthaul interface to split the Next Generation Node-B (gNB) into i) a remote unit implemented using a Universal Service Radio Peripheral (USRP) from ETTUS research (subsidiary of National Instruments), ii) a Digital Unit (DU) and iii) a Central Unit (CU) using OAI-5G [8]. Note that OAI is an open-source project that implements 3GPP compliant-technology on general purpose x86 computing hardware and a USRP. Initially, OAI implemented the 4G architecture, but OAI wide community have started implementing the 5G new radio, the NF-approach and the control-user plane separation for the 5G Core (5GC). Recall that Docker [5] is a tool that enables developers to create, deploy and run applications by using containers. Containers package up an application with all its dependencies (libraries and environment) and ship it out as only one package making a vast economy of scale. Unlike Virtual Machines (VMs), containers virtualize an application on top of the operating system kernel directly, and thus containers are lighter to

load, faster to start and less memory-hungry, when booting, compared to VMs.

B. Control Layer (SDN Controllers)

An orchestrator organizes containers at the networking level and enables the application to run as intended. We used Docker Swarm [5] that is the native Docker orchestrator as it is well-suited for small-scale deployments. In addition, it has a large community attached to it due to the huge number of developers that use Docker products. As Controllers, we use FlexRAN [11] as a Software-Defined Radio Access Network Controller (SD-RANC) to manage the gNB and Open Network Operating System (ONOS) as software-defined transport network controller to manage the fronthaul, midhaul and backhaul networks.

C. Intelligent application Layer (Top)

On first hand, we have developed a northbound web app interfacing with the controllers for slices lifecycle management. It orchestrates the configuration management tasks, including on-demand slices creation/update/deletion, cell re-configuration, devices rehomeing to slices and performance statistics generation. This layer enables us to create on-the-fly different slices as needed (eMBB/mMTC/uRLLC). Based on our previous work where we exhaustively evaluated multiple ML techniques [3], we have shortlisted a particular regression model, as we will see in section V-B, and integrated it to automate the provisioning of network slicing ratios. Furthermore, we have used BeC3 [4] for IoT applications to manage and control mutiple IoT Devices. BeC3 offers an intuitive interface that hides the complexity and specificities of IoT objects, while giving unified access to their core functionality. Using this, we abstract the peculiarities of IoT devices in a virtual representation to get a reusable VM or a Docker image of such IoT devices. BeC3 allows implementing universal typical commands for IoT devices (such as ON/OFF for smart switches, GET-TEMP for temperature sensor, GET-LUMINOSITY for light sensors, and so on). It also offers the possibility to control the behavior of integrated IoT devices, such as programming a device to upload measurement data to a remote server.

D. TICK Stack (Vertical)

In order to manage the performance of the virtualized network functions in real-time, all of Telegraf, InfluxDB, Chronograf and Kapacitor (TICK) [3] are used as an end-to-end performance stack. Telegraf collects time-series data from a variety of sources and in particular, from FlexRAN API through JSON structured data describing the radio access network metrics. InfluxDB delivers high performance read/write access and efficiently stores fetched time-series data. Chronograf is used to visualize graphs of stored data in InfluxDB using different plots formats. Finally, we used Kapacitor to send commands to the northbound application. TICK stack provides extract-transform-load capability from heterogeneous raw data and allows autonomous anomalies detection in time-series data.

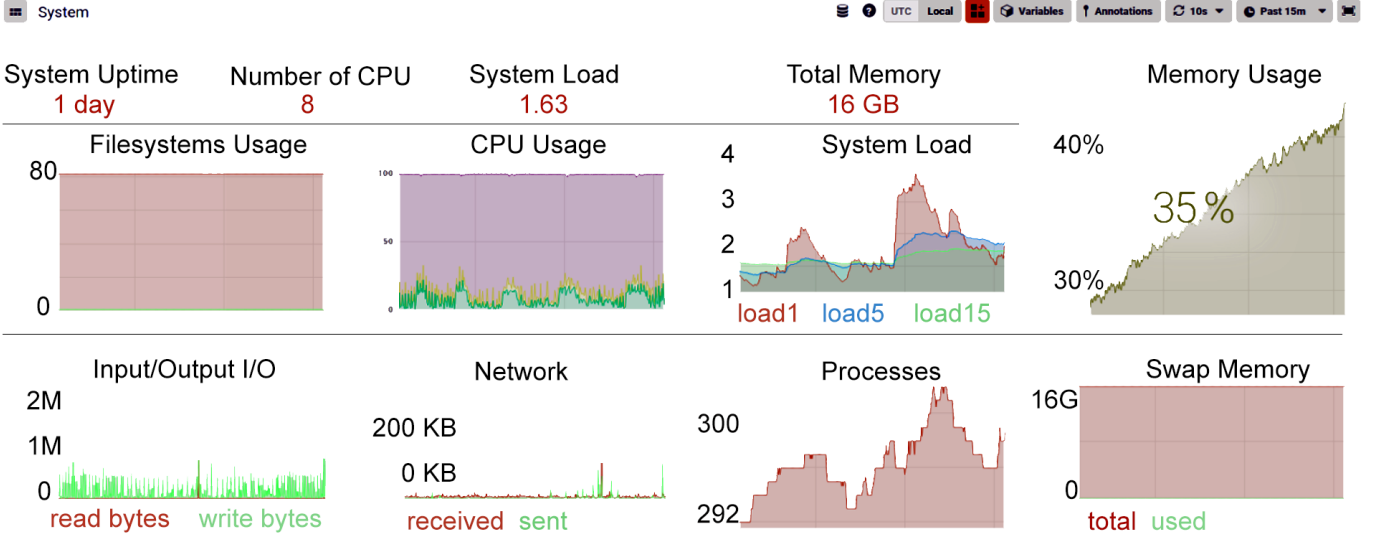


Fig. 2. Physical Machine Metrics Evolution over Time

Fig. 2 depicts our hardware layer performance dashboard including system uptime, number of cores of the physical machine, system load, total memory, hard disk activity in terms of Input/Output (I/O) read/write, network activity, processes and swap. All these metrics depict our physical host system capabilities (8 Central Processing Units (CPUs), 16 GB of memory, etc), when implementing our different uses-cases, shown in the next section. We can notice that the memory usage increases over time and reads 35% of the available memory due to the caching process.

Note that this prototype is depicting a Non Stand-Alone deployment of 5G, but it could be easily upgraded to a standalone one or any other form of the eight envisioned deployment options of 5G. It is worth noting that our 5G experimental prototype takes into consideration all of: self-healing, high-availability, centralized coordination, and horizontal scalability. Indeed, the swarm Cluster for 5GC, provides self-healing capability by re-instantiating a replica of a microservice whenever the latter fails due to intermittent software glitch or software/hardware failure of a worker machine and so on. To do so, it automatically recreates another microservice to compensate for the failed one and maintain the same requested number of microservices of the service in question. Using a cluster composed of a manager and multiple workers ensures high-availability of the 5GC, meaning that if a worker fails, the cluster is maintained and remaining workers and manager absorb and balance the load, using the internal or an external load-balancer such as HAPROXY [15] (our case). The latter plays two roles of load-balancing and reverse-proxying. This service-based architecture ensures seamless scalability by means of horizontal scaling, precisely, by scaling-out or scaling-in the number of replicas according to changing needs in terms of computing resources. Note that the details of our auto-scaling algorithm can be found in our previous work [15].

V. PERFORMANCE EVALUATION

A. Dynamic Slicing and Devices Auto-Rehoming

In this work, we have focused on both eMBB and mMTC types of slices. Indeed, since wireless bandwidth is scarce and limited, eMBB slices are anticipated to be scaled in chunks on two fronts: bandwidth and service functions. On the other hand, massive sensors typically collect small-data payloads and exchange them with a base station. Due to the limited power of embedded sensors envisioned for mMTC slices, it is preferred that they work intermittently rather than continuously. Accordingly, the advantages of a dynamic slicing of radio resources are vivid in these two types of 5G network slices, as we will see in the following.

It is worth noting that our proposed architecture could also address the uRLLC network slices since our RAN network functions are deployed as microservices in a container-based virtual environment. According to our experiments [13], we have quantified the average deployment time of a container-based function to $1.8 \mu s \pm 0.2 \mu s$, which makes our prototype able to meet the critical time constraints of uRLLC.

In this context, we have used our 5G experimental prototype [3] to manage several IoT Devices including smart switches, smart lamps and temperature/light sensors. We used our BeC3 IoT platform to provision the IoT devices at both virtual and physical levels. Knowing that BeC3 is lightweight, we deployed it onto a small-device (Raspberry Pi 2 model B) to act as an IoT Gateway connected to the 5G prototype either using Wi-fi and a tethered mobile internet from a connected smartphone or through a USB dongle equipped with a corresponding SIM card.

We depict multiple use cases in Fig. 3. In particular, use case (A-i) illustrates the signal flow of a virtual wattmeter configured to send periodic measurements to BeC3 cloud, where the results are displayed using a virtual display, both provisioned using BeC3. Use case (B-i) illustrates a manually triggered event using a smart-switch to control a smart lamp

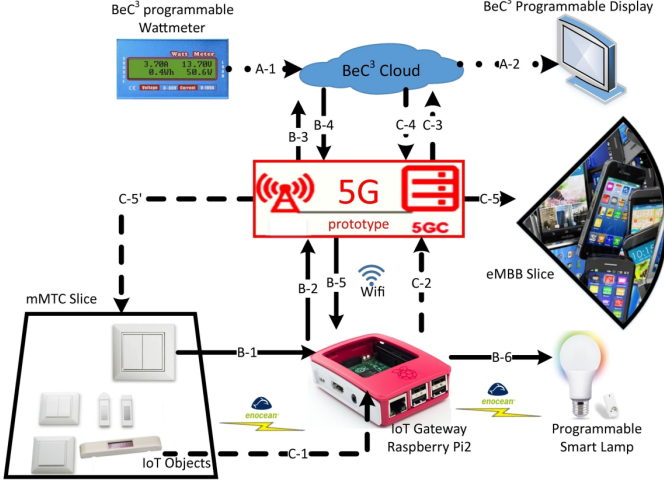


Fig. 3. Use cases for evaluation of our prototype

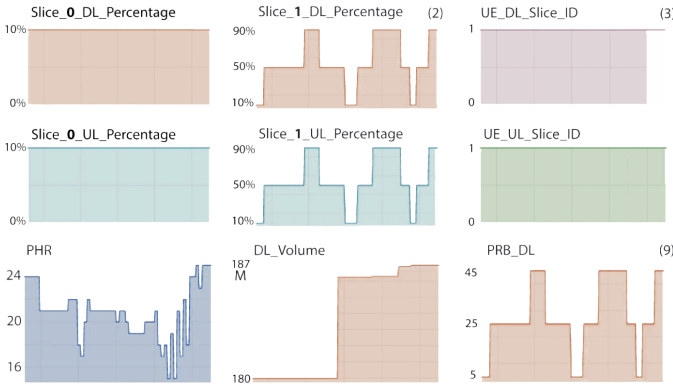


Fig. 4. Radio Metrics Dashboard over Time

through mobile internet provided by our 5G experimental prototype. The smart devices use EnOcean serial protocol 3 to transmit related radio telegrams. Finally, use case (C-i) shows an automated process of dynamic slice creation and a rehoming of IoT devices to an mMTC slice, while parenting smartphones to an eMBB slice.

The results of the dynamic slicing in Chronograph dashboard can be seen in Fig. 4. For the sake of clarity, we chose a step of 40% of the available Physical Resource Blocks (PRBs).

As the dynamicity of throughput requirements is more interesting in the case of the eMBB slice than that for the mMTC slice, we considered configuring dynamic auto-scaling of PRBs on the eMBB slice. We can see, in placeholder (2), that Slice 1, configured to use dynamic slicing, exploits the auto-scaling by increasing its ratio from 10% to 50% then 90% of available PRBs and vice versa. We can also see the corresponding allocated PRBs are (5, 25 and 45) out of the available 50 PRBs of our setup, as reflected in placeholder (9), which depicts the instant downlink PRBs. Remaining placeholders depict the downlink/uplink percentages for both of slice 0 and 1, in addition to the instant Power-HeadRoom (PHR) and the downlink volume which increases along the usage.

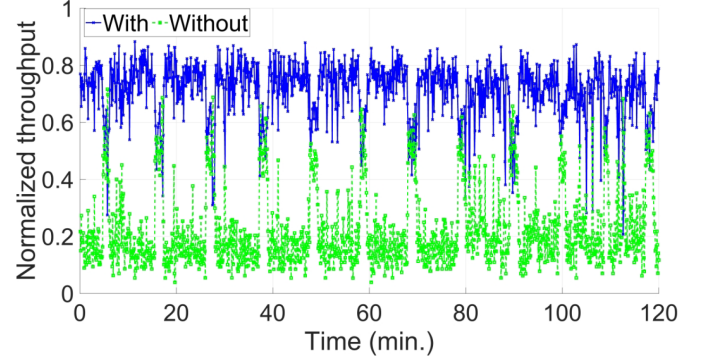


Fig. 5. Normalized Throughput with/without slicing ratios forecasting

B. Slicing Ratio Forecasting using Machine Learning

Using our proposed prototype, we initially acquired considerable RAN data to train multiple ML models. In particular, we have selected regression trees as our preferred ML forecasting mechanism as we have found it to be the best performer in terms of Root Mean Square Error (RMSE), Mean Average Error (MAE) and coefficient of determination (R^2) [3].

Recall that a regression tree is built using a binary recursive partitioning, which is an iterative process that splits the data into partitions or branches. Then, it continues splitting each partition into smaller groups as the mechanism moves up through each branch.

Using our prototype, and precisely, by carrying-out a query to the database (InfluxDB) tables, we extracted the historical traffic profile including the timestamp, and the used PRBs by each connected device. We then aggregated the used PRBs per type of slice. We have also enriched such traffic profile by appending the requirements from the slice owner. All of these constituted our set of predictor variables. On another hand, according to the implemented infrastructure usage policy (relying on iterative auto-scaling), we have appended to our enriched traffic profile the related response for each slice in terms of allocated PRBs. Accordingly, we have fed this whole table to the regression tree learner.

We have implemented, in the northbound web application, our ML based forecasting mechanism to orchestrate slicing ratios for multiple network slices based on available PRBs. Once the training of our ML model is done, our objective would be to dynamically provision optimum-slicing ratio out of the available pool of PRBs to the new created slices. To do so, we start by creating one slice hosting a smartphone and map it to an eMBB slice. We use Youtube video streaming application to simulate some bursty traffic. Note that this slice is configured to use our dynamic slicing through forecasting, whereas a second slice is dynamically created, whenever the IoT devices have to upload some data streams. This latter slice will be then automatically deleted after finishing the transfer process and its related PRBs will be returned to the PRBs pool. Accordingly, the remaining active slices can use them when needed. We have evaluated the system performance using two metrics, namely the normalized throughput and the CPU utilization of the host machine.

Fig. 5 depicts the normalized throughput for the eMBB Slice

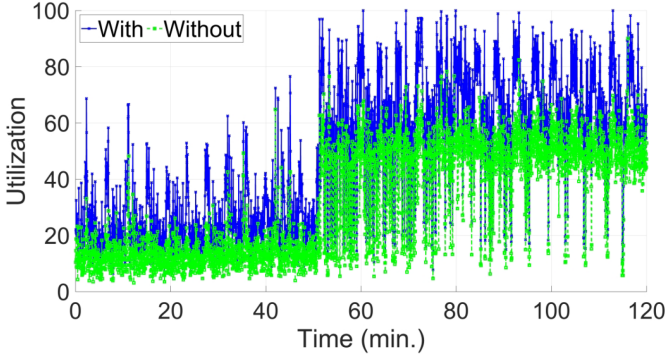


Fig. 6. CPU utilization with/without slicing ratios forecasting

during two hours of simulation when using our autonomous network slicing versus the case where we deactivate it. The normalized throughput is computed as the ratio between the obtained throughput and the maximum available one. We can easily observe that the forecasting process allowed increasing the throughput by approximately 30% compared to the case where the forecasting is de-activated and only a static slicing ratio of 50% to each of these two slices is used. We notice that forecasting slicing ratios accelerates the provisioning of required PRBs to serve the requirements of the running applications.

On the other hand, we can see that this automation comes at a cost of an increased utilization of CPU of the host system as depicted in Fig. 6. This is explained by the increase of the computing resources utilization, resulting from the processing of the slice lifecycle (creation/destruction of a slice). Finally, we note that the same normalized throughput is maintained, although we queued a new video, starting from 50 minutes and onwards. This is seen in the system utilization, where the trends of CPU usage for both cases (with/without slicing ratios forecasting) increased, reflecting such additional load.

VI. CONCLUSION

In this article, we investigated the available prototypes for 5G and proposed a novel architecture. Our architecture differs from existing state-of-the-art monolithic ones as it relies on a microservice-based implementation. We also complemented it with a TICK Stack for performance monitoring and management. On application layer, we developed two web applications for configuration management, leveraging machine learning forecasting techniques to auto-provision slicing ratios. We proposed to manage the IoT devices using BeC3. Using Docker Swarm, we implemented self-healing, high availability, centralized coordination, and horizontal scaling mechanism. We validated our implementation using a realistic use-case with IoT sensors, so that we dynamically create a slice and rehome devices to it. We observed that, using a machine-learning forecasting model, has substantially increased the throughput of the network, based on same radio network design, with a drawback of an increased computing resources utilization.

ACKNOWLEDGMENT

This work was supported by FUI SCORPION project (Grant no. 17/00464), Paris Ile-de-France Region (Grant no. 19CNE00437), “Azim & Saade Foundation”, and the Lebanese University.

REFERENCES

- [1] A. A. Barakabitze, A. Ahmad, R. Mijumbi, and A. Hines, “5G network slicing using SDN and NFV: A survey of taxonomy, architectures and future challenges,” *Computer Networks*, vol. 167, p. 106984, 2020.
- [2] D. S. Linthicum, “Cloud-native applications and cloud migration: The good, the bad, and the points between,” *IEEE Cloud Computing*, vol. 4, no. 5, pp. 12–14, 2017.
- [3] N. Salhab, R. Rahim, R. Langar, and R. Boutaba, “Machine learning based resource orchestration for 5g network slices,” in *2019 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2019, pp. 1–6.
- [4] S. Cherrier, I. Salhi, Y. M. Ghamri-Doudane, S. Lohier, and P. Valembois, “BeC3: behaviour crowd centric composition for IoT applications,” *Mobile Networks and Applications*, vol. 19, no. 1, pp. 18–32, 2014.
- [5] “Docker,” <https://www.docker.com/>, accessed: 2021-01-10.
- [6] E. Coronado, S. N. Khan, and R. Riggio, “5G-EmPOWER: A software-defined networking platform for 5G radio access networks,” *IEEE Transactions on Network and Service Management*, vol. 16, no. 2, pp. 715–728, 2019.
- [7] C.-Y. Huang, C.-Y. Ho, N. Nikaein, and R.-G. Cheng, “Design and Prototype of A Virtualized 5G Infrastructure Supporting Network Slicing,” in *2018 IEEE 23rd International Conference on Digital Signal Processing (DSP)*. IEEE, 2018, pp. 1–5.
- [8] “OpenAirInterface,” <https://openairinterface.org/>, accessed: 2021-01-10.
- [9] “O-RAN,” <https://www.o-ran.org/>, accessed: 2021-01-10.
- [10] “MOSAIC-5G,” <https://mosaic5g.io/>, accessed: 2021-01-10.
- [11] X. Foukas, N. Nikaein, M. M. Kassem, M. K. Marina, and K. Kontovasilis, “Flexran: A flexible and programmable platform for software-defined radio access networks,” in *Proceedings of the 12th International Conference on emerging Networking EXperiments and Technologies*, 2016, pp. 427–441.
- [12] X. Foukas, F. Sardis, F. Foster, M. K. Marina, M. A. Lema, and M. Dohler, “Experience building a prototype 5G Testbed,” in *Proceedings of the Workshop on Experimentation and Measurements in 5G*, 2018, pp. 13–18.
- [13] S. Matoussi, I. Fajjari, S. Costanzo, N. Aitsaadi, and R. Langar, “5G RAN: Functional split orchestration optimization,” *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 7, pp. 1448–1463, 2020.
- [14] “Kubernetes,” <https://kubernetes.io/>, accessed: 2021-01-10.
- [15] N. Salhab, R. Rahim, and R. Langar, “NFV Orchestration Platform for 5G over On-the-fly provisioned Infrastructure,” in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, 2019, pp. 971–972.

Nazih Salhab (nazih.salhab@univ-eiffel.fr) is a subject matter expert on mobile networks working in Amaris Consulting. He received his dual Ph.D. from University Paris-Est, France and the Lebanese University. Before that, he received his Master 2 and Engineer’s degree from the Lebanese University. He has been leading business analysis, operation and project management for mobile network operators since more than 15 years. His research interests are 5G/6G, network slicing, C-RAN, SDN, NFV and machine learning.

Rami Langar (rami.langar@univ-eiffel.fr) is currently a Full Professor at University Gustave Eiffel (UGE), France. He received his Ph.D. degree in network and computer science from Telecom ParisTech, Paris, France in 2006. His research interests include resource management, network slicing in 5G/6G systems, C-RAN, SDN, machine-learning, and mobile Cloud offloading.

Rana Rahim (rana.rahim@ul.edu.lb) is an Associate Professor at the Lebanese University. She received her Ph.D. degree in January 2008 from the University of Technology of Troyes (UTT) - France. Her research interests include 5G/6G, QoS, IoT, Smart Grids, C-RAN, network slicing and SDN.

Sylvain Cherrier (sylvain.cherrier@univ-eiffel.fr) is an Associate Professor at University Gustave Eiffel (UGE), France. He received his Ph.D. degree from University Paris-Est, France. His research interests include Internet of Things IoT management, choreography and automation.

Abdelkader Outtagarts (abdelkader.outtagarts@nokia-bell-labs.com) is a senior researcher and leader in Nokia Bell Labs. He received the M.Sc. and Ph.D. degrees in automation of energy systems from INSA Lyon, France, and the M.Sc. degree in software engineering from the ETS of Montreal. His research interests include 5G/6G, NFV, SDN, and machine learning.