



Intelligent Orchestration of End-to-End Network Slices for the Allocation of Mission Critical Services over NFV Architectures

Bego Blanco¹(✉) , Rubén Solozabal¹ , Aitor Sanchoyerto¹ ,
Javier López-Cuadrado¹ , Elisa Jimeno² , and Miguel Catalan-Cid³

¹ University of the Basque Country, Bilbao, Spain
{begona.blanco,ruben.solozabal,aitor.sanchoyerto,javilo}@ehu.eus

² Atos, Madrid, Spain
elisa.jimeno@atos.net

³ i2CAT Foundation, Barcelona, Spain
miguel.catalan@i2cat.net

Abstract. The challenge of deploying mission critical services upon virtualised shared network models is the allocation of both radio and cloud resources to the critical actors who require prioritized and high-quality services. This paper describes the design and deployment of an intelligent orchestration cycle to manage end-to-end slices on a NFV architecture. This novel tool includes the monitoring of the network elements at different levels and the processing of the gathered data to produce the corresponding alert mitigation actions.

Keywords: Network slicing · Orchestration · NFV · Mission-critical

1 Introduction and Related Work

5G networks are expected to bring a new disrupting ecosystem, prompting the creation of innovative next generation vertical applications. To that end, one of the most awaited features is the provisioning and management of network slices tailored to the needs of each particular vertical industry and specific deployment. In particular, Network Function Virtualization (NFV) is embraced as one of the key technologies that will allow the creation of customized network slices to meet different service requirements.

The public safety sector will be one of the major beneficiaries of this technological development. Traditional mission critical applications expose tight QoS requirements, which find difficulties to be fulfilled by traditional network models. In consequence, traditional public safety networks have demanded private and dedicated network models, which eventually lead to an inefficient use of resources

and spectrum. But now, network slicing through proper resource orchestration is making the network sharing model a reality.

The concept of network slicing was introduced by the Next Generation Mobile Network (NGMN) alliance within its whitepaper [8]. Later, the 3GPP took the responsibility of standardising this technology, defining the entities and the functionality required to manage network slicing [2]. Presently, network slicing is integrated in the ETSI-NFV architecture [4]. Current NFV standards [5] define the interaction between the network slice management functions defined by the 3GPP and the NFV Management and Orchestration (MANO) module, establishing the required connection between the network controllers and the NFV orchestration unit in order to perform the dynamic assignment of network resources.

However, the implementation of the concept of orchestrating a service slice within this standardized network architecture is still in a development phase. In this sense, there are some independent initiatives as [3, 6, 7, 9] that are contributing to the creation of modules that complement the current MANO capabilities in order to orchestrate E2E slices.

In this paper, we present an NFV-based intelligent orchestration cycle with the capability of providing a set of shared resources to deal with the dynamic reconfiguration challenge. This orchestration cycle has been developed in the scope of H2020 5G ESSENCE project [1]. The slice concept introduced in 5G, along with the highly virtualised and software-based deployments, enables the automatic on-the-fly adjustment of the resource assignment to the changeable environment. This feature is of utmost importance in mission critical applications where sudden events can instantly alter the network requirements and priorities. For this reason, this work provides a comprehensive approach to demonstrating dynamic End-to-End (E2E) slices reconfiguration and service adaptation in a mission critical deployment.

The paper is organised as follows: Sect. 2 describes the orchestration cycle defined to dynamically adjust the end-to-end network slices in a NFV-based deployment. Next, Sect. 3 describes the validation scenario to later discuss the obtained results. Finally, Sect. 4 summarizes the main contributions and poses new research challenges that will be addressed in the future.

2 NFV-Based Intelligent Orchestration Cycle

NFV comes up driven by the telecommunications industry in order to enhance the deployment flexibility, foster the integration of new services within operators and also attain CAPEX/OPEX drawdowns.

However, the dynamic allocation of resources to separated and customised network slices still remains a challenge. This section describes a novel orchestration cycle providing new tools for automated E2E network slicing. The proposed orchestration cycle involves the monitoring system, the alert mitigation module and the execution of the mitigation actions.

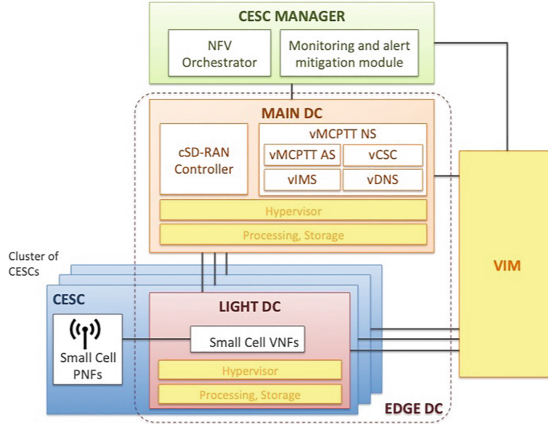


Fig. 1. 5G ESSENCE network architecture.

2.1 Network Architecture

The 5G ESSENCE approach, depicted in Fig. 1, takes the existing 5G architectures as a reference point, combining the 3GPP framework for network management in Radio Access Network (RAN) sharing scenarios and the ETSI NFV framework for managing virtualised network functions. Our architecture allows multiple network operators (tenants) to provide services to their users through a set of Cloud-enabled Small Cells (CESCs) deployed, owned and managed by a third party (i.e., the CESC provider). The CESC offers virtualised environment with computing, storage and radio resources at the edge of the mobile network. This cloud can also be ‘sliced’ to enable multi-tenancy.

Besides, the two-tier architecture of 5G ESSENCE is well aligned with the 5G architecture described by 5G-PPP, where the infrastructure programmability is identified as one key design paradigm for 5G. First, 5G ESSENCE achieves infrastructure programmability by leveraging the virtualised computation resources available in an Edge Datacenter (Edge DC). These resources are used for hosting VNFs tailored according to the needs of each tenant, on a per-slice basis. Second, the Main Datacenter (Main DC) allows centralising and softwarising the control plane of small cell functions to enable a more efficient utilisation of the radio resources coordinated among multiple CESCs.

We propose to enhance the orchestration functionalities adding more intelligence into the CESC Manager (CESCM) together with the NFV Orchestrator (NFVO). In particular, 5G ESSENCE provides a network monitoring and alert mitigation mechanism that supports and improves both the NFVO and RAN controlling functions. The event flow for the management of end-to-end slicing for a Mission-critical Push-to-talk (MCPTT) service is depicted in Fig. 2, and each component is further described in the following sections.

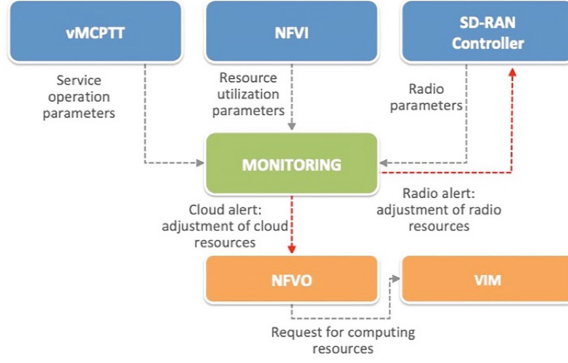


Fig. 2. End-to-end slicing event flow.

2.2 System and Service Monitoring and Alert Mitigation

The main objective of the Monitoring and Alert Mitigation system shown in Fig. 3 is to access the available information about the network elements and process it in order to conclude if and when a network reconfiguration is needed.

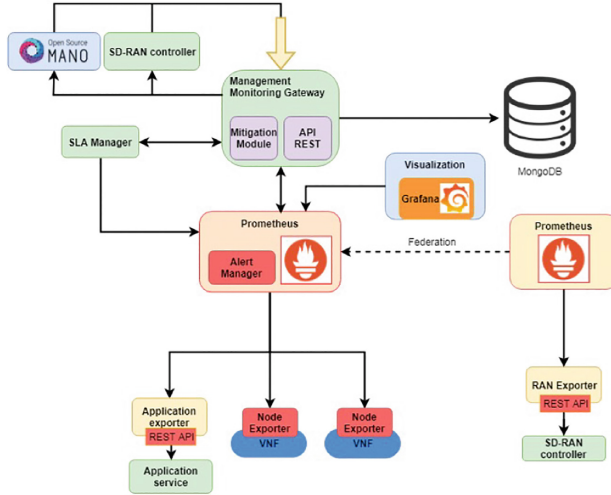


Fig. 3. Monitoring and alert mitigation architecture.

The orchestration cycle begins with the collection of the monitored data through the **exporters** in each monitoring-enabled building block. The monitored data is stored in Prometheus, which is on charge of triggering the alerts as defined according to the different services and their Service Level Agreement

(SLA). These alerts are defined to notify about an unexpected behaviour in the system and SLA violations.

It must be also noted that the monitoring of the Wi-Fi RAN controller relies on the **federation** of the Prometheus server installed in the component. Federation allows Prometheus to have a heritage of some targets monitored from another Prometheus. The main idea for using Federation is to have a decentralised system in order to monitor the Wi-Fi RAN metrics through another Prometheus for other tasks.

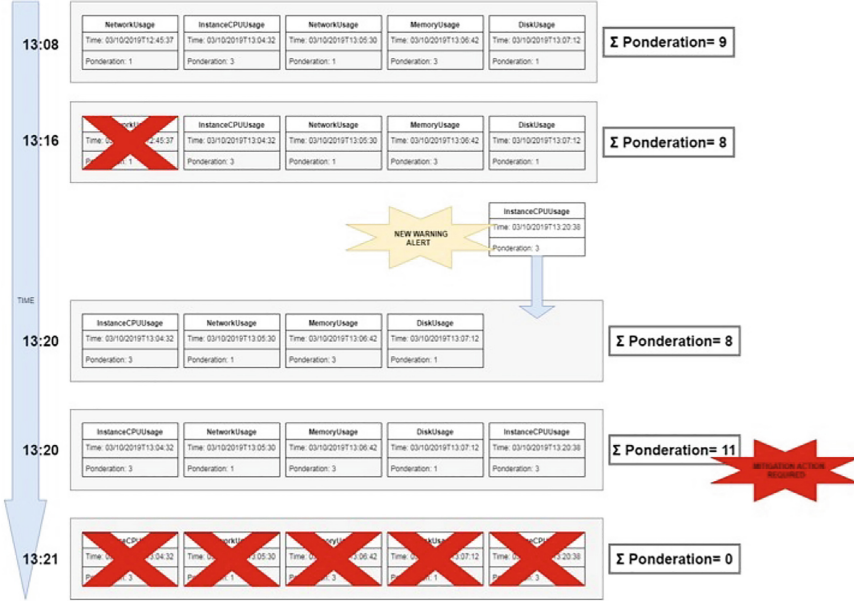


Fig. 4. Flow for mitigation warning alerts.

The alerts raised are picked by the **Alert Mitigation Module (AMM)**, which is part of the **Management Monitoring Gateway**. The purpose of AMM is to manage the configuration of the architectural components responsible of the behavior of the E2E slice. To that aim, AMM contains the mitigation logic based on a ponderation of the rules defined in the **Rulebook**. When an alert is triggered, AMM differences between different severity levels. If the severity is critical, the mitigation module must mitigate the alert with higher priority without considering further alerts following the configuration defined in the Rulebook. For warning severities, the Mitigation module saves the alert in a time window, which is configured by the Rulebook (Fig. 4). The window (or queue), groups the alerts by the specific mitigation required by it. Every warning alert has a ponderation in the Rulebook. The warning alert is added in the queue with its correspondent ponderation. If the sum of the ponderations

in the mitigation queue exceeds the mitigation ponderation, configured in the Rulebook, a mitigation action composed with all the warning severity alerts is triggered, emptying the mitigation queue and silencing the alert.

Finally, within the scope of this paper, we have defined two endpoints to forward the mitigation actions and close the monitoring and mitigation loop: the NFVO orchestrator to manage the scaling options of the Network Service, and Wi-Fi RAN controller to manage the resources used by the Wi-Fi slice. These two blocks are further described in the next sections.

2.3 MCPTT Service Architecture

Mission Critical Push-To-Talk (MCPTT) is a mission critical communication standard that allows half duplex one-to-many and one-to-one voice services in order to coordinate emergency teams. Users request permission to transmit pressing a button. Additionally, the MCPTT service provides a means for a user with higher priority (e.g., MCPTT Emergency condition) to override (interrupt) the current speaker. As it appears, the management of this type of half-duplex communication is not trivial, since it requires an appropriate management of priorities and privileges to allow communication.

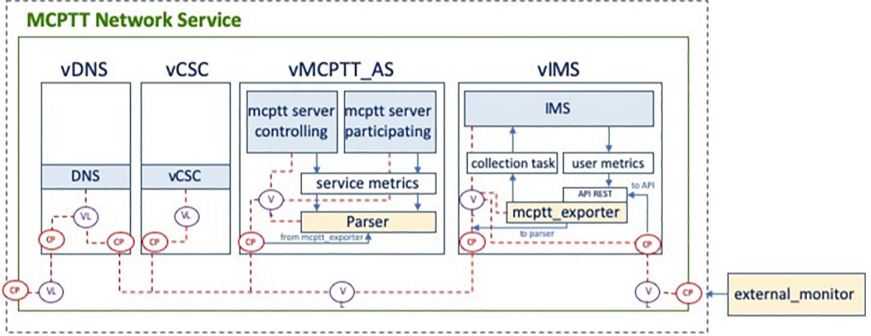


Fig. 5. MCPTT service architecture.

The MCPTT Network Service is composed of one VNF that completes the mission critical push to talk service. This service is defined in multiple Virtual Deployment Units (VDU) to optimise the usage of the resources: a DNS server, an IMS (IP Multimedia Subsystem) service for session management, a CSC (Common Service Core) for service status information, and the MCPTT AS (Application Server) providing centralised support for MCPTT services and call control. Figure 5 depicts the deployment of the described MCPTT network architecture.

In order to integrate the described MCPTT network service within the orchestration cycle detailed above, we must include a tailored exporter to extract

the required metrics for the monitoring tasks. It appears as *mcptt_exporter* in Fig. 5. This component is responsible for collecting the metrics from the MCPTT service to later expose them for the analysis in the monitoring system. It is implemented as a REST API: when the *mcptt_exporter* receives a status request from the Prometheus in the Monitoring module, it queries the involved components of the NS (mainly IMS and *MCPTT_AS*) to gather the metrics and format them properly.

2.4 RAN Controller Architecture

Figure 6 depicts the components of the RACOON Wi-Fi RAN slicing solution.

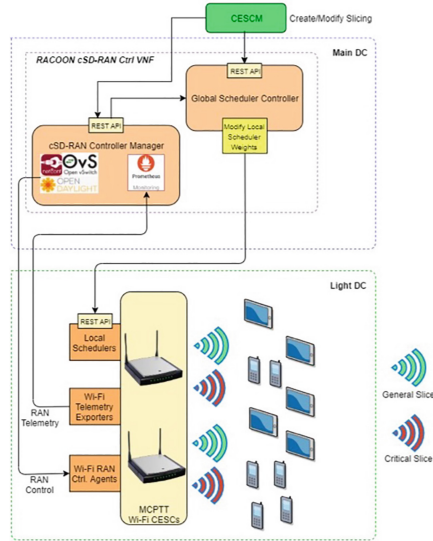


Fig. 6. RACOON SD-RAN controller architecture.

The **Controller Manager** is the core of RACOON. It is in charge of OpenDayLight SDN controller, Open vSwitch database server and the Netconf Manager by means of the different implemented clients (REST APIs) and controls the CESCes according to the deployed slices and services. It also gathers telemetry from the Wi-Fi RAN by means of its Prometheus server. Moreover, through its REST API, the RAN Controller exposes the management of the infrastructure and the slices to the CESCM.

The **Global Scheduler Controller** manages the weights/quotas of the instantiated slices in the Wi-Fi RAN. It allows enabling, modifying and disabling the local schedulers of the different Wi-Fi CESCes, which locally manage the percentage of airtime or channel time assigned to each slice. It implements a REST API to allow its control via the CESCM.

Finally, the **Wi-Fi CESC** is composed by Single Board Computers (SBCs) with a Linux distribution. The main software used in order to deploy Wi-Fi connectivity is Hostapd, which has been modified in order to deploy, monitor and control multiple virtual Access Points (vAPs) on top of a single physical interface, according to the desired Wi-Fi slices. By means of these modifications, the Local Scheduler is able to manage the MAC-scheduler which controls the airtime or channel time assigned to each slice (which is then fairly distributed among all the user terminals of each slice). Also, it hosts a Prometheus Exporter (Hostapd Exporter¹) in order to gather RAN telemetry.

3 Orchestration Cycle for MCPTT Deployments

This section shows the results of the integration of the described enhanced orchestration tools developed within the 5G ESSENCE project to deploy a MCPTT service slice. To that aim, we first declare the metrics collected from the network elements and the mitigation actions defined when an emergency event is detected. Then, we describe the validation scenario and show the results of the complete deployment.

3.1 Monitoring Metrics and Mitigation Action Definition

The monitoring system collects network status information from network elements at different levels: NFVI through *node_exporter*, MCPTT service through a tailored exporter and RACOON cSD-RAN controller through Prometheus federation. The information collected from the NFVI that is involved in this experiment includes CPU, memory and disk usages, VM port throughput and availability of VMs (if they are up). The information collected from the MCPTT service includes the number of registered and active users, the number of private calls, group calls, private emergency calls, group emergency calls that have been started/ongoing/terminated and the number of users involved in each of the calls. Finally, the system collects information from RACOON about the number of users per slice per cell, transmitted bit rate and quality of the signal. For each identified metric, the measurement framework and the alarms it can trigger is included. Two alarm thresholds are defined. The first alarm threshold provides a warning, whereas the second threshold is considered a critical situation.

3.2 Scenario Definition and Deployment Results

The demonstration of the dynamically orchestrated MCPTT deployment cannot be based on a static scenario, since one of its objectives to be proven is the elastic allocation of resources attending to different levels of emergency conditions detected by the monitoring system. We propose a deployment topology in three main stages.

¹ http://bitbucket.i2cat.net/users/miguel_catalan/repos/hostapd_prometheus_exporter/browse.

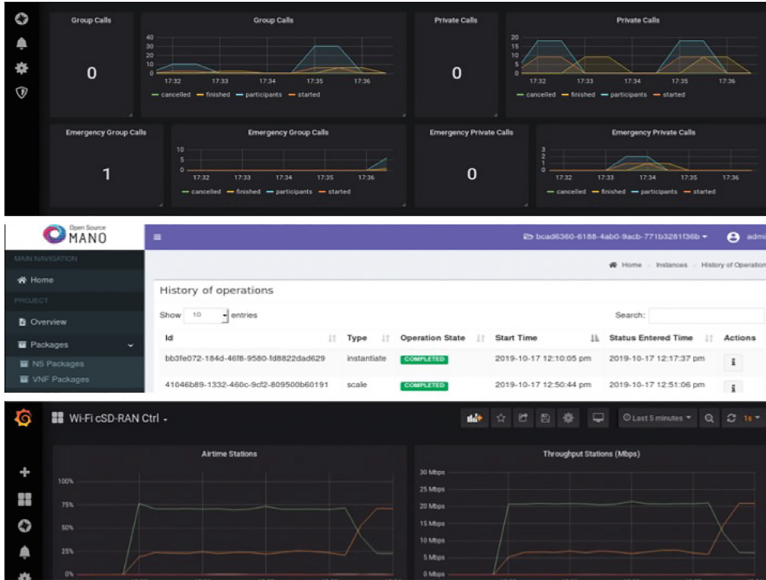


Fig. 7. Orchestration cycle screenshots.

Under normal circumstances, the system instantiates the network slices that correspond to a default service agreement. Here, the first responder only needs a reduced amount of access capacity and communication features for its normal operations. Then, triggered by an emergency incident that is detected through a private emergency call, the first responder requires increased capacity in terms of edge computing resources, in order to serve a higher number of incoming communications and/or public safety users. This implies the scaling of the MCPTT VNF and it may involve a deterioration of the service for legacy users, since their network slice(s) must be reduced in order to appropriately allocate the higher priority MCPTT service. Finally, in the third stage triggered by a group emergency call, the system responds with an expansion of the MCPTT radio slice up to the 75% of the available bandwidth in the cell where the emergency events are happening (detected by the increasing number of users in the cell). Again, this situation may involve an impairment of the service provided to civilians in favor of the communications for first responders, which require higher priority.

Figure 7 shows some screenshots that illustrate the operation of the orchestration process. The upper screen shows the monitoring of MCPTT calls during the experiment. It can be observed how the different events are detected over time. The screenshot in the middle shows the result of an alert mitigation action in the second stage that leads to the MCPTT VNF scale. Finally, the lower screenshot shows the reconfiguration of the radio slice as a result of the mitigation action in the third stage.

4 Conclusions and Future Work

This paper has described the intelligent orchestration cycle that proves that the 5G ESSENCE context provides a solution for an efficient and elastic E2E network slicing and the efficient orchestration of the radio and cloud resources. The results highlight the value of the shared network model, demonstrating the capacity of the 5G ESSENCE architecture to autonomously allocate resources to first responders whenever they are required, but giving them up to the commercial services when the requirements are low. The elastic allocation of resources is performed automatically, leveraging the monitoring and alert mitigation functionalities that complement the orchestration processes in the CЕСSM.

Our research work will continue to further develop orchestration tools to enhance the E2E slicing capabilities of NFV environments. New research trends include the use of machine learning techniques in the decision-making process, the migration and placement of VNFs and the analysis of the possibilities of multi-RAT access.

Acknowledgement. This work has been partly funded by the EU funded H2020 5G-PPP project 5G ESSENCE (Grant Agreement No 761592).

References

1. H2020 5G ESSENCE project. <https://www.5g-essence-h2020.eu>
2. 3GPP: Study on management and orchestration of network slicing for next generation network. TR 28.801, 3GPP, January 2018
3. Chien, H.T., Lin, Y.D., Lai, C.L., Wang, C.T.: End-to-end slicing as a service with computing and communication resource allocation for multi-tenant 5G systems. *IEEE Wirel. Commun.* **26**(5), 104–112 (2019)
4. ETSI: Network Function Virtualization (NFV). Architectural Framework. GR GS NFV 002, ETSI, October 2013
5. ETSI: Report on Network Slicing Support with ETSI NFV Architecture Framework. GR NFV-EVE 012, ETSI, December 2017
6. Khalili, H., et al.: Network slicing-aware NFV orchestration for 5G service platforms. In: 2019 European Conference on Networks and Communications (EuCNC), pp. 25–30. IEEE (2019)
7. Montero, R., Agraz, F., Pagès, A., Spadaro, S.: End-to-end network slicing in support of latency-sensitive 5G services. In: Tzanakaki, A., et al. (eds.) ONDM 2019. LNCS, vol. 11616, pp. 51–61. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-38085-4_5
8. NGMN: Description of network slicing concept (2016)
9. Ni, R., et al.: An end-to-end demonstration for 5G network slicing. In: 2019 IEEE 89th Vehicular Technology Conference (VTC2019-Spring), pp. 1–5. IEEE (2019)