



HAL
open science

Quality of perception prediction in 5G slices for e-Health services using user-perceived QoS

Yosra Benslimen, Joanna Balcerzak, Albert Pagès, Fernando Agraz, Salvatore Spadaro, Konstantinos Koutsopoulos, Mustafa Al-Bado, Thuy Truong, Pietro G Giardina, Giacomo Bernini

► **To cite this version:**

Yosra Benslimen, Joanna Balcerzak, Albert Pagès, Fernando Agraz, Salvatore Spadaro, et al.. Quality of perception prediction in 5G slices for e-Health services using user-perceived QoS. 2020. hal-02572024

HAL Id: hal-02572024

<https://hal.science/hal-02572024>

Preprint submitted on 13 May 2020

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.

Quality of perception prediction in 5G slices for e-Health services using user-perceived QoS

Yosra Ben Slimen⁽¹⁾, Joanna Balcerzak⁽¹⁾, Albert Pagès⁽²⁾, Fernando Agraz⁽²⁾, Salvatore Spadaro⁽²⁾, Konstantinos Koutsopoulos⁽³⁾, Mustafa Al-Bado⁽⁴⁾, Thuy Truong⁽⁴⁾, Pietro G. Giardina⁽⁵⁾, Giacomo Bernini⁽⁵⁾

⁽¹⁾ *Orange Labs, Chatillon, France*

⁽²⁾ *Universitat Politècnica de Catalunya (UPC), Barcelona, Spain*

⁽³⁾ *Creative Systems Engineering, Athens, Greece*

⁽⁴⁾ *Dell EMC, Cork, Ireland*

⁽⁵⁾ *Nextworks, Pisa, Italy*



Quality of perception prediction in 5G slices for e-Health services using user-perceived QoS

Yosra Ben Slimen⁽¹⁾, Joanna Balcerzak⁽¹⁾, Albert Pagès⁽²⁾, Fernando Agraz⁽²⁾, Salvatore Spadaro⁽²⁾, Konstantinos Koutsopoulos⁽³⁾, Mustafa Al-Bado⁽⁴⁾, Thuy Truong⁽⁴⁾, Pietro G. Giardina⁽⁵⁾, Giacomo Bernini⁽⁵⁾

⁽¹⁾ *Orange Labs, Chatillon, France*

⁽²⁾ *Universitat Politècnica de Catalunya (UPC), Barcelona, Spain*

⁽³⁾ *Creative Systems Engineering, Athens, Greece*

⁽⁴⁾ *Dell EMC, Cork, Ireland*

⁽⁵⁾ *Nextworks, Pisa, Italy*

Abstract

In order to compete for a prominent market share, network operators and service providers should retain and increase the verticals' subscription, catering to their needs in order to differentiate themselves from competitors. In this scenario, verticals' satisfaction arises of paramount importance. As such, user experience is becoming a reliable indicator for service providers and telecommunication operators to convey overall end-to-end system functioning. To properly estimate end user satisfaction, operators and service providers require efficient means for quality monitoring and estimation at all layers, in conjunction with mechanisms able to maintain said quality at optimum levels. Given these factors, this paper proposes a mechanism for Quality of Perception (QoP) estimation in e-Health services, enabling the QoP-aware management of network slices fulfilling the requirements of supported services. To this end, the paper proposes a cognitive-based architecture which allows for the collection and monitoring of verticals' data to estimate QoP and provides mechanisms to re-configure the underlying network slices according to the monitored quality levels. A machine learning (ML) model is introduced that aims to forecast any future degradation in the quality perceived by vertical users. In case of a predicted degradation, the proposed architecture reacts and triggers the necessary remedial actions, referred as actuations. In order to evaluate the developed ML model and to showcase the

interaction between the different components of the proposed architecture, an experimental study is presented with real data extracted from a roaming ambulance. In addition, a Proof of Concept of the actuation mechanism is demonstrated through an experimental testbed emulating e-Health services.

Keywords: 5G network slicing, quality of experience, machine learning, network cognitive management.

1. Introduction

5G subscriptions are predicted to reach 1.9 billion by 2025 where 35% of traffic will be carried by 5G networks [1]. The novelty of 5G technology is the concept of network slicing that should allow an increased speed, performance, scalability, and flexible service deployment. According to [2], a network slice (NS) consists of physical and/or virtual network functions (PNF/VNF) that can belong to access and core network part. Then, this network functions are interconnected by means of network resources (what is known as chaining), composing a synthetic infrastructure with specific characteristic, both in functionalities and resource capacities. The synthesis of a NS, then, serves a particular functional purpose and once instantiated, it is used to support certain communication services, which ultimately are deployed to support vertical services on top.

E-Health is identified as one of the key vertical sectors to which 5G infrastructures should focus in providing solutions for its effective materialization. In particular, 5G could provide support for stroke diagnosis through the provisioning of dedicated NSs. This would help a trained paramedic team to assist stroke patients with the cooperation of medical personnel (e.g. a doctor) present in remote facilities, such as hospitals.

To date, significant progress has been made on in-hospital stroke management, but a reliable pre-hospital in-ambulance solution has not yet been established. Hence, solving this problem could offer life-saving time gains and speed up treatment initiation by early activation of the in-hospital stroke response, thereby curtailing the risk of misdiagnosis and death. Generalised in-ambulance tele-medicine is a recently developed and promising approach, which is under the umbrella of the more generic e-Health vertical industry.

Following 5G service classification, e-Health services belong to the Enhanced Mobile Broadband (eMMB) type of services defined by 5G standards. These type of services require the provisioning of significant network capac-

ities in a highly mobile and distributed geographical area. E-Health has as added requirement strict Quality of Experience (QoE) service levels, which ultimately may have an impact on the well-being of stroke patients.

The targeted scenario in this paper begins with the continuous collection and streaming of patient data when the emergency ambulance paramedics arrive at the incident scene. Several Internet of Things (IoT) devices (e.g. wearables, cameras) enable the provision of enhanced patient insights. The goal is for all paramedics to have wearable clothing that can provide real-time video feeds as well as other sensor related to data pertaining to the immediate environment. Figure 1 summarizes the scenario in hand.

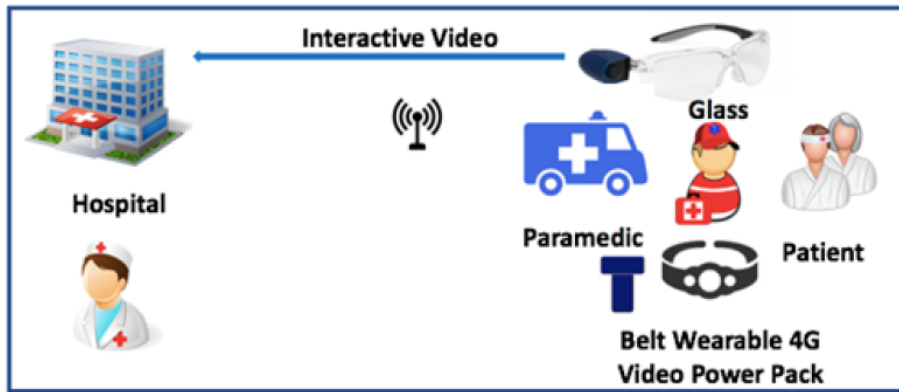


Figure 1: Summary of e-Health use case.

Under such scenario, the network should provide efficient communication capabilities that enable storing and potential real-time streaming of patient data to the medical team inside the ambulance. This can enable more intelligent decision support for the paramedics attending the patient. Therefore, detecting and resolving e-Health NSs’ operational problems is of a paramount importance in order to minimise the downtime of slice operations and to provide a good QoE.

Quality of Experience has been defined by [3] as the degree of delight or annoyance of the user of an application or service. Nevertheless, a good quality of service does not guarantee the satisfaction of the customers. As defined in [4], the Quality of Service (QoS) is the characteristics of a telecommunications service that bears on its ability to satisfy stated and implied needs of the user of the service (e.g., packet loss, latency, etc.). Using these two definitions, it is clear that the customer experience matters the most in QoE.

In the context of network slicing, measuring customer experience (i.e. slice user) is a challenging problem for verticals because it is costly and complex due to the human involvement in the process. The challenging issue is mainly the subjectivity of this measurement. First, applications are of different importance to different users and they also have diversified network-level QoS requirements. Second, different application-level QoS performances bring different effects on user QoE. Third, many contributing factors of QoE change over time so they are highly time-variant. Therefore, it is inadvisable and far from reality to measure the QoE once forever.

In addition to that, although a good QoS does not imply a good QoE, the impact of network QoS parameters, such as delay and packet loss rate are considered as the contributing factors that influence the user perceived QoE. For example, user satisfaction can be expressed by reducing the traffic delay that can cause dropping of sessions. Thus, the network QoS parameters (i. e. bandwidth, delay, packet loss rate, etc.) are regulating the user perceived quality experience.

In this paper, we utilize a more objective user-centric approach that computes the User-perceived QoS also called ‘the Quality of Perception’ (QoP). We aim to provide an architecture that provides a top QoP to the vertical. A machine learning (ML) model is presented in order to analyse the vertical’s QoP and to forecast the degradation of perceived signal strength quality that may arise in the future 5 minutes. In case of any predicted degradation, an alert is sent to an actuation framework that will react accordingly in order to correct the problem before it occurs.

The remaining of this paper is structured as follows: Section 2 reviews the main work found in the literature in regards of cognitive management of NSs and quality estimation in network services. Section 3 introduces the proposed architecture. Section 4 elaborates the e-Health use case and the data collection mechanism. Section 5 details the supervised forecasting ML model for QoP degradations. Section 6 introduces the actuation framework. Section 7 presents the testbed, along the numerical experiments that aim to evaluate the performance of the ML model as well as the details of the actuation mechanism in order to correct potential degradation of the measured slice QoP. Finally, Section 8 elaborates the conclusions and future works.

2. Related work

Network slicing is a very hot topic in the research and industrial communities, due to the many benefits that it brings. In this regard, the cognitive management of network slices is seen as a promising candidate to the automated and autonomous fulfillment of vertical services requirements, such as e-Health. Thus, several works can be found tackling the multiple aspects of slice management in 5G shared infrastructures. In this section, we review some of the work found in the literature regarding network slicing in 5G, cognitive network management and quality estimation for vertical services.

2.1. Network slicing in 5G

Network slicing in the context of 5G is a concept that has been introduced by the Next Generation Mobile Network (NGMN) Alliance [5]. Network slicing consists on the partitioning of the underlying physical infrastructure, network elements (both physical and virtual) as well as functions to enable the deployment of self-contained networks, each one tailored to the specific needs of the services that will run on top.

Such a concept enables multi-tenancy in multi-domain/role/technology 5G infrastructures, in which all involved infrastructure segments are able to create isolated network slices towards the fulfillment of end-to-end (E2E) synthetic infrastructures, each one managed and controlled independently from the dynamic of the underlying substrate.

As such, several works has tackled efficient ways to facilitate slicing in the multiple segments that constitute the 5G ecosystem. Starting from the Radio Access Network (RAN) segment, works have focused on the challenges posed by the wireless nature of the RAN medium, in order to effectively slice the RAN spectrum and provide slices with the requested resources. For instance, authors in [6] focus on the resource assignment aspect of RAN slices, proposing allocation algorithms that coordinate both the resource assignment to slice instances as well as the sharing across provisioned slices.

Another critical segment is the Mobile Edge Computing (MEC), since it allows the deployment and execution of network functions close to the endpoints, reducing drastically the latency for communication services that are delay sensitive (e.g. e-Health). Thus, it is essential to devise architectures that facilitate the slicing and allocation of functions in MEC sites. To this end, authors in [7] present an architecture specially tailored to the challenges posed by network slicing in MEC segments, with the aim of minimizing

communication latencies and signalling overheads for IoT services. These works, among others focused in other segments (e.g. Core, optical transport, multi-domain), provide the necessary foundations that enable the realization of 5G slices.

2.2. Cognitive network management

Aside from the aspects related to enable slicing in 5G networks from an infrastructure perspective, another important topic relates to the management of said slices. Due to the sheer scale of 5G systems, the dynamicity, heterogeneity and volatility of the supported traffic profiles, as well as the added complexity of softwarized/virtualized infrastructures, traditional management solutions are no longer suitable for the management of network slices. As such, novel management solutions should be capable to automatically and autonomously adapt to changes or unforeseen states on the underlying infrastructure and supported services.

In this regard, cognition-based management, focused on Artificial Intelligence (AI)/ML procedures, is gaining momentum as a solution to overcome the challenges posed by network slice management in 5G environment. The importance of investigating AI/ML-based network management for 5G networks has been recognized by several large European Union projects under the 5G PPP in recent years. For instance, in the Phase 1 of the 5G PPP program, the SelfNet project [8] proposes a Software Defined Networking (SDN)-/Network Function Virtualization (NFV)-based network management framework for advanced Self-Organizing Network (SON) capabilities in 5G infrastructures. Similarly, another 5G PPP Phase 1 project CogNet [9] also targets AI/ML-based network management solutions and provides a service portfolio including data gathering and pre-processing, quality assurance, network demand prediction, location-based services and planning.

In addition, there are several works in literature that target different facets of cognitive management in softwarized network environments. Focusing on the aspect of media delivery in 5G networks, authors in [10] propose a system that leverages on ML for a network resource allocator system as the main contribution which enables autonomous network management with quality-awareness. As a summary of the application of ML/AI techniques in SDN/NFV environments, authors in [11] review the challenges and opportunities of ML/AI in softwarized network environments, with a special focus on data-driven decision making for management and control of SDN/NFV-based infrastructures. To this end, the authors propose to enhance the func-

tional primitives of monitoring, composition and control with ML modules. All these works, among many others, pave the path towards cognitive-aided management of networks, specially in the presence of SDN/NFV technologies.

2.3. Quality of Experience/Service estimation

Due to the more user oriented nature of 5G networks, a concept that has been named as "verticals in the loop", quality assurance becomes primordial for an optimal delivery of services towards the vertical customers. While network level QoS is quite easy to monitor/estimate, with a plethora of mechanism and strategies to guaranty good QoS levels, the quality perceived by end users is more challenging to determine and optimize, since it involves the subjective perception of human actors about the performance of a delivered services.

As such, several works in literature propose strategies in support of quality estimation and optimization at all levels, ranging from systems that aid on decision taking for optimal service provisioning, to QoE estimation in media/network services, passing through works related to network-level quality monitoring and estimation.

For instance, authors in [12] propose a novel NFV orchestration solution that employs ML techniques to enhance the Quality of Decision (QoD) of an NFV-based system when orchestrating and provisioning services. The objective is to enhance NFV systems to achieve a near-optimal placement of virtualized network functions at minimum monitoring costs.

Focused on the estimation of QoS in cloud environments, authors in [13] review the main open questions about QoS attributes modeling and prediction in cloud services and showcase the limitations posed by statistics-based techniques in their forecasting. In this regard, authors discuss how ML-based techniques arise as a superior alternative in the presence of dynamic workloads for proper QoS estimation.

In the aspect of QoE estimation, the main line of research tackled in the literature is QoE estimation/prediction in video delivery services. Indeed, video delivery is a primordial service in nowadays telecom/cloud ecosystem and economic activity, encompassing both consumer grade services (e.g. video on demand, video chat) or business grade services, as it would be the case of e-Health. Thus, several works explore the main challenges of this aspect and propose several solutions to overcome them. The main common denominator is the use of ML/cognitive-aided solution, since the own nature

of QoE requires of mechanisms that are able to learn what is the perceived quality from a user perspective. As an example, the work presented in [14] employs ML in order to enhance traditional RAN bases stations with QoE awareness, improving the delivery of video streams in current and future Long Term Evolution (LTE) radio networks. In general, QoE management for video delivery services is an E2E process, which requires not only QoS management at network level but also real time monitoring and control of end user QoE. Therefore, traditional strategies relying on statistical or reactive measures fall short in addressing properly the posed challenges. In this regard, the work presented in [15] reviews the most significant QoE management methods in video services, pinpointing to the most effective ML methods employed in the literature in this regard.

3. Overview of the cognitive management architecture

In this section, we present a management architecture which aims on enabling 5G systems with cognitive-based management of NSs in support of advanced e-Health vertical services, with quality guarantees assurance being the main pillar in which the architecture is built. The presented architecture rests within the framework of the 5GPPP Phase II SliceNet project [16]. SliceNet aims to provision NSs with QoE guarantees to support vertical industries. As such, it proposes a management/control architecture devoted to QoE-aware life-cycle management of 5G slices. Figure 2 depicts a schematic of the architecture, putting emphasis on the functionalities and roles that are detailed in the remaining of this section.

The depicted architecture considers a multi-role environment, in which several stakeholders with clear responsibilities engage for the delivery of 5G NSs. The following sub-sections provide more details about the key components and entities involved.

3.1. Network Service Providers

A network service provider (NSP) is the owner of the physical infrastructure. Multiple NSPs operate under a "network slice as a service" (NSaaS) business model, offering their infrastructure capabilities to an upper layer (i.e the Digital Service Provider). In this regard, the NSPs are responsible for the maintenance of the physical network infrastructure, which may encompass several segments, for example, RAN, MEC, core and NFV Points of Presence (such as data-centres).

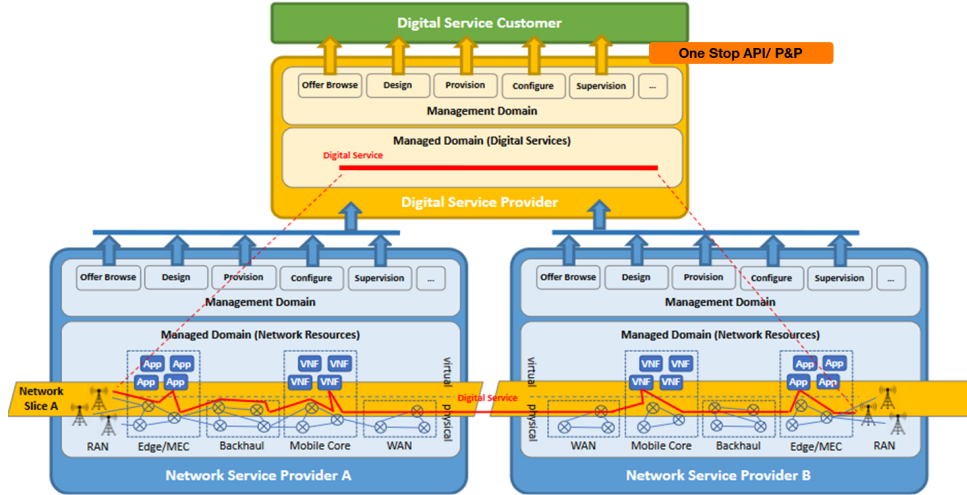


Figure 2: Overview of SliceNet architecture and stakeholders relationships.

In this regard, NSs deployed at NSP level serve the networking requirements that are needed within a single domain towards the composition of E2E services, i.e. a NS offered to vertical customers. To this end, the NSPs have in place control and management elements that enable the provisioning, configuration and monitoring of physical/virtual resources that constitute NSs. These NSs are bound by specific Service Level Agreements (SLAs) that guarantee the properties of the deployed slices towards the fulfilment of specific needs requested by upper layer entities.

3.2. Digital Service Providers

A Digital Service Provider (DSP) is responsible for the creation and management of E2E NSs to be offered in a NSaaS model towards vertical customers to satisfy the needs for supporting their associated services. Therefore, SliceNet architecture divides its functionalities between the multiple roles (NSP/DSP). While single-domain control functionalities are basically assumed within the NSP role, the E2E multi-domain management functionalities are integrated within the DSP role.

This latter aims to translate the vertical service requirements towards specific network/infrastructure requirements, which then are translated onto a specific E2E NS to be fulfilled. The E2E NSs are then broken down into single domain NS which are materialized thanks to the capabilities offered by the multiple NSPs that have a business relationship with the DSP.

3.3. Access and FCAPS management

SliceNet adopts a vertical (i.e. slice user) centric approach that aims at abstracting technology and domain related enablers towards the provision of services tailored to the needs of the topmost user of the infrastructure. In this respect a Plug and Play (P&P) framework has been designed with the aim of provisioning vertical oriented exposure of the delivered slice. This requires that specially crafted control plugs are developed and available to be selected by verticals when requesting the provisioning of their slices. Each of the control plugs provides a specific customized northbound slice control view that fits the needs of the vertical. A dedicated slice P&P control instance (composed by one or more plugs) is automatically created as part of the slice provisioning process, and then exposed to the vertical. Each P&P control instance is able to process, at its southbound, the platform specific information as this is maintained via the layered and multi-domain Fault, Configuration, Accounting, Performance and Security (FCAPS) enablers. The FCAPS framework is also provided in an abstraction fashion starting from the pillar technologies per NSP domain and spanning up to inter-domain aggregations as administered by the DSP domains. The overall governance of the several domain enablers and abstractions is applied through the One Stop Application Programming Interface (OSA) that spans vertically from NSP to vertical. Besides, OSA provides role based views of the artifacts that have to be administered under a producer consumer relationship between overlaid roles.

In practice FCAPS allows NSPs to abstract their domain specific technologies with respect to actuation and sensing capabilities and expose these as offerings through the OSA towards DSPs. DSPs can consume the monitoring information as this is arranged per slice via the NSPs Data Lakes as well as invoke actuation information through the layered orchestration entities. DSPs maintain higher level slice counters and metrics, in their own Data Lakes, as these have been made available through exposure functions instantiated per slice. At the level of DSP, Data Lake information is subject to be processed by Cognitive, QoE and P&P functionalities that have been activated per slice to deliver higher level and novel FCAPS strategies.

Beyond the bottom-up approach for the delivery of slice specific metrics and due to the vertical centric approach, verticals can contribute to the Data Lake information for the exploitation of service metrics that the user endpoints are able to produce. In this respect the vertical tailored user interface (UI) view provided by P&P and OSA is augmented to create slice specific

endpoints that expose the existing P&P options towards more automated (UI less) procedures to contribute to the slice segment of the DSP Data Lake and trigger cognitive and QoE functions via direct quantitative vertical feedback. OSA allows therefore slice owners to indicate which P&P functions should be exposed as Representational State Transfer (REST) endpoints that can be invoked by metric and counter agents deployed in vertical user equipments (UEs) and collecting valuable metrics from the user domain.

3.4. Cognition sub-plane

It is precisely in the DSP role in which SliceNet defines a novel Cognition Sub-Plane, incorporating AI/ML capabilities, for the QoE-aware management of E2E NSs. The role of the Cognition Sub-Plane is dual. On one hand, it monitors the multiple metrics and events that are relevant to the quality of deployed slices. Given this information, it derives elaborated data that allows to predict the quality of slices thanks to specialized ML models tailored to the characteristics of NSs and vertical services. On the other hand, given the alerts generated by the models, the Cognition Sub-Plane is responsible to apply remedial actions so as to maintain optimal quality levels. This is done through a policy-driven actuation framework, which requests for specific (re-)configurations to selected NSP. It aims to modify the characteristics of the underlying single-domain NSs in order to have an optimal E2E NS able to provide the topmost quality levels towards vertical customers.

4. e-Health vertical use case

Following the architecture previously depicted in Figure 2, the E2E slice customer is a national/regional health service organization (in our case National Ambulance Service in Ireland), which operates multiple (static) hospitals, dispatch centres, and (moving) ambulances. The E2E e-Health slice offered to the customer by the DSP consists initially of a “base” NS Instance (NSI) containing the minimal set of network functions and services. The base slice is fairly static and is centered around a geographical area in the vicinity of the hospital. The hospital/dispatch hosts experts who provide real-time support to the paramedics.

As ambulances are dispatched, additional Network Sub-Slice Instances (NSSIs) may be instantiated in order to increase the geographical coverage of the slice or to guarantee the latency and availability requirements of the slice. For the latter, additional processing functions may be dynamically

instantiated at suitable MEC locations. Dispatch may trigger a handover of a paramedic’s communication stream to a different hospital. Handover between domains might be needed while the ambulance is moving. The vertical customers only see services provided by their DSP.

For data collection, we installed G-Net Pro application in a Samsung phone and configured the application to capture the statistics for down-link traffic within 1s. Following the ambulance paths, the application collected the network statistic information including the timestamp, longitude, latitude, mobility speed, technology (e.g., 4G, 3G, WiFi), technology mode, bit-rate, cellID, serving cell, operator ID, etc. in multiple files. We then applied a data pre-processing program to extract the interesting fields for data training which is described in Table 1. In total, we have 12 data-sets, each data-set contains 2-3 hours statistic information.

Table 1: Collected QoP metrics from the ambulance

Feature	Description
Perceived Reference Signal Receive Power (RSRP)	Average power received from a reference signal
Perceived Reference Signal Receive Quality (RSRQ)	It indicates the quality of the received signal
Perceived Signal to Noise Ratio (SNR)	It indicates the signal strength relative to background noise
Perceived Channel Quality Indicator (CQI)	It indicates how good/bad the communication channel quality is
Perceived Received Signal Strength Indicator (RSSI)	It indicates the signal strength measured from all base stations
Perceived DL-bit-rate	Bit-rate of the download link
Perceived UL-bit-rate	Bit-rate of the upload link

As shown in Figure 3, the information collected from the UEs at the ambulance are fed towards the monitoring systems at the NSPs, since UEs are connected to a specific NS within the E2E NS, for which a concrete NSP is responsible. The collected information is then transferred to the DSP, more specifically, to the Anomaly Detection ML model that resides in the Analyser.

This model acts as an intelligent sensor that observes the last 5 minutes of the QoP metrics of the NS supporting the e-Health services. The sensor

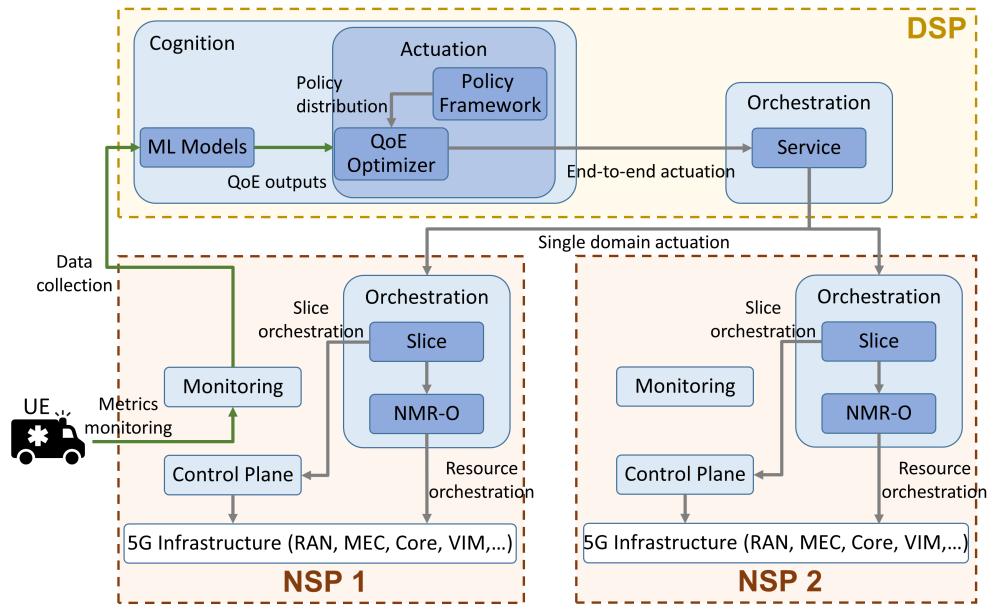


Figure 3: SliceNet e-Health use case actuation workflow.

will predict every 5 seconds if a degradation in the network signal strength in the RAN segment may be perceived by the vertical in the future 5 minutes as depicted in Figure 4. Based on the alerts generated by the model, the actuation framework reacts and triggers the necessary remedial actions.

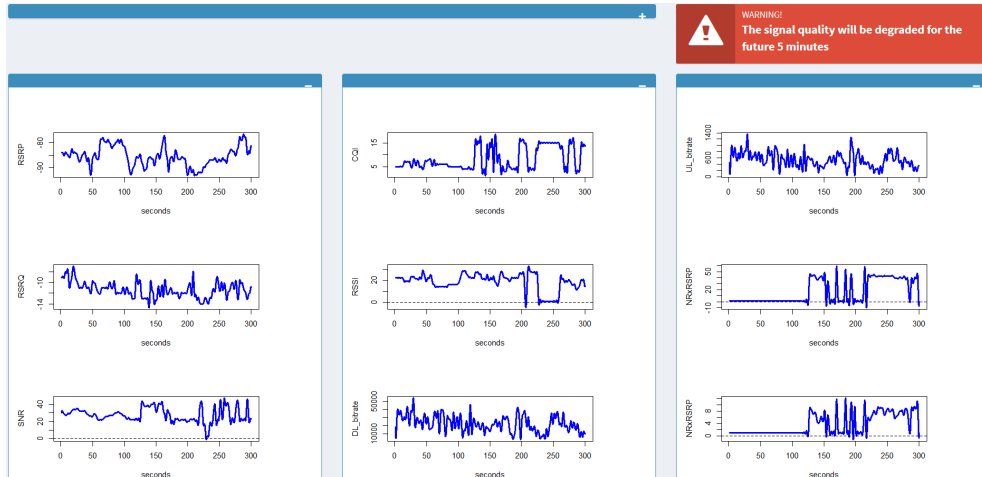


Figure 4: QoP prediction model objective

Due to the mobility nature of ambulances, which host the UE end-points of the NS, RAN connectivity is of paramount importance, since it allows to connect the medical equipment to the NS regardless of the location as long as it is under the coverage of the provisioned E2E NS. A bad quality of signal will disturb the services running on top of the slice since it may cause a disruption of the communication between the paramedics, a degradation of the vertical quality of experience and a danger for the patient. As such, if a bad signal quality is predicted, the ML model needs to alert the the actuation framework that will trigger remedial actions.

Such a loop enables the real-time processing of the monitored UE data, its subsequent analysis and near-real-time actuation to overcome the predicted communication degradation. These requirements are essential in the e-health scenarios since these latter offer critical services that need a fast and up-to-date analysis and actuation to guarantee optimal quality levels.

Given the overall scenario and goal, the following sections proceeds on detailing the elaborated QoP prediction model and the actuation framework.

5. Quality of Perception prediction model

The data \mathbf{X} under study are a set of n observations as illustrated in Figure 5. Each observation X_i corresponds to the curves of QoP metrics observed for the last 5 minutes. Hence, each observation is described by a set of p curves and a label: $\mathbf{X} = \{(\mathbf{X}_i(t), label_i)\}_{t \in [0, T], 1 \leq i \leq n}$ with $\mathbf{X}_i(t) = (X_{i1}(t), \dots, X_{ip}(t))' \in \mathbb{R}^p$, $p \geq 1$ and $label_i = 1$ if there is a signal degradation will be perceived by the vertical in the prediction horizon (the future 5 minutes) given $\mathbf{X}_i(t)$, 0 otherwise.

5.1. Imbalance management

Due to the imbalance in the class, the learning of the different patterns that are hidden in the data may be misled and it must be treated before the training phase. For this reason, we use the SMOTE technique [17] that is a well-known approach for over-sampling the minority class in a data set. The over-sampling is achieved by generating "intelligent" copies of the minority class observations i.e. by artificially creating synthetic samples. The new examples are generated by using the nearest neighbors of these observations. In consists in perturbing one attribute at a time by a random amount within the difference to the neighboring instances. This approach effectively forces the decision region of the minority class to become more general.







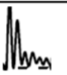





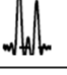

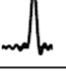



KPI 1	KPI 2	...	KPI p	Label
		...		1
		...		0
...
		...		1
...
		...		0
		...		1
...
		...		0

Figure 5: An illustration of the training set structure

5.2. Functional data extraction

In real life, the QoP metrics are extracted with discrete time points (every second). Therefore, a smoothing [18] is applied to every set of QoP metrics observed for the last 5 minutes in order to reflect their functional nature in a continuous time interval. Smoothing assumes that each observed curve x_{ij} ($1 \leq i \leq n$, $1 \leq j \leq p$) can be expressed as a linear combination of basis functions $\{\phi_l\}_{l=1,\dots,M}$: $x_{ij}(t) = \sum_{l=1}^M a_{ijl}\phi_l(t)$, $t \in [0, T]$, where $\{a_{ijl}\}_{l=1,\dots,M}$ are the basis expansion coefficients. As a result, every KPI curve is a functional object that is represented by M values. These latter are the coefficients of the corresponding QoP curve after its projection in the new functional space.

5.3. Dimension reduction

By considering that the numbers of the observed QoPs and the observation duration might be huge, a dimensionality reduction seems inevitable.

Therefore, a Functional Principal Components Analysis for multivariate data (FPCA, [18]) is applied. It consists in computing the principal components C^h and principal factors f^h of the Karhunen-Loeve expansion: $X(t) = \mu(t) + \sum_{h \geq 1} C^h f^h(t)$, $t \in [0, T]$. It allows to have an optimal representation of curves from the set of functional data into a functional space of reduced dimension. The number of principal components $m \leq M$ is chosen so that at least 90% of the information is covered. As a result, each curve is defined by a vector of its principal components of size m . For more details, we refer the reader to [18] and [19].

5.4. Supervised learning

At this stage, the QoP measurements are transformed in terms of principal components, a training data with numeric values and a binary class is ready for a training phase. The training is assured by a Random Forest [20] which is a classification technique that does not require much memory. Random forest run-times are quite fast and it is able to deal with imbalanced and missing data. It starts by splitting the data set into a number of samples. For each sample, a decision tree is trained. It is a "divide-and-conquer" approach that aims to form a "strong learner" (a forest) from a group of "weak learners" (decision trees). From the set of the decision tree classifiers that have been trained on various sub-samples of the data set, the random forest uses averaging to improve the predictive accuracy and to control over-fitting. To classify a new observation, it puts this latter down each of the trees in the forest. Each tree gives a classification i.e. it "votes" for a class and the forest chooses the classification having the most votes.

In case of imbalanced data sets, the prediction error between classes is highly imbalanced. One class have a low prediction error compared to the other. Then random forests, by trying to minimize overall error rate, will keep the error rate low on the large class while letting the smaller class to have a larger error rate. As a result, the accuracy of the model will not be misled by the imbalance. Figure 6 illustrates the training phase of the model.

Given the smoothing basis, the FPCA basis and the learned random forest, prediction of future anomalies for new observations is possible as described in Figure 7.

For each new observation, a smoothing is first applied for every QoP metric. The observation in terms of functional features is then projected on the same FPCA basis used in the training phase. Given the obtained principal components, the learned random forest will then predict whether a

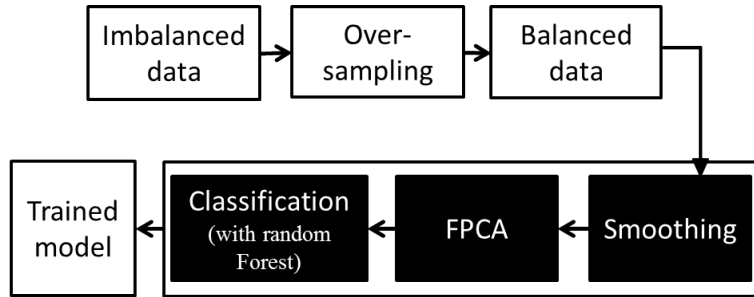


Figure 6: Training phase

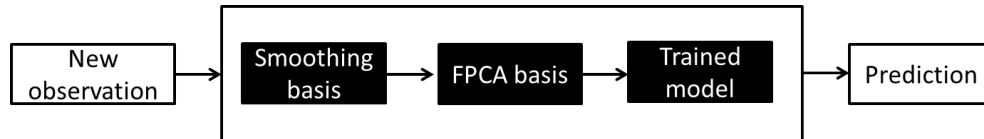


Figure 7: Prediction phase

degradation of the signal quality will be perceived in the future 5 minutes or not.

6. Actuation framework for QoP optimization

Aside from the presented model, which opens the door for QoP-aware management of e-Health NS, it is also important how the model engages with other elements in order to close the cognition loop as to maintain top notch quality levels. With this in mind, an actuation framework has been developed within the presented architecture inside the cognition sub-plane.

Following the explanations of previous sections, and to complete the cognitive loop for QoE-aware management of e-Health NSs, this section elaborates about the developed actuation framework to be triggered as response to the outputs from the presented QoP prediction ML model. As illustrated in Figure 3, the actuation framework consists in a Policy Framework and a QoE Optimizer, which serves to execute remedial actions once degradation on the slices quality is observed.

The QoE Optimizer is the element responsible to trigger remedial action in cases which the quality of deployed E2E slices is compromised or is not within satisfactory levels. To do so, it listens to events (outputs) of the ML models within the cognition sub-plane (for the case at hand, the developed

QoP prediction model) and, in case the events mark that quality levels of the slice are not as expected or some anomaly can affect them, it triggers the necessary (re-)configurations of the E2E NS and composing NSSIs to regain optimal quality levels.

To do so, the QoE Optimizer interacts with the service orchestrator at the DSP level (Service), which, in turn interact with the multiple single-domain slice orchestrators present at the multiple NSPs. The slice orchestrators may engage directly with control plane functions if configuration details of the NSs need to be modified (e.g. change the assigned bandwidth, apply a UE handover) or may engage with the resource orchestrator at the same NSP (NMR-O) in case that characteristics of the virtual resources need to be modified (e.g. CPU/memory assigned to a virtual function). Thanks to that, the full actuation is orchestrated.

This presented mechanism is governed through policies in place at the QoE Optimizer. These policies are disseminated from a Policy Framework at the DSP level, following an Event-Condition-Action (ECA) model [21] which specifies for what event (an output from the ML model) and condition (e.g. an anomaly is predicted) what action should be applied (e.g. increase the bandwidth at the RAN segment). Thanks to that, it is possible to follow a rather declarative approach in which the Actuation Framework specifies the "what" and "when" of the remedial action while the "how" is implemented thanks to the logic of the elements responsible for orchestration or configuration of slices/infrastructure elements (i.e. orchestrator, control plane).

Aside from that, policies also carry information about the end-points that need to be listened for monitoring the information related to the events as well as a small logic that dictates simple extraction/aggregation operations for this information, named as Sensors. In this way, all the elements that belong to the monitoring pipeline will be properly configured following the specifications of the Sensors present in the policies for a specific E2E NS.

7. Numerical experiments and testing

In the previous sections, we presented the architecture of the QoP-aware management for e-Health 5G NS, the ML model in charge of QoP predictions and the Actuation Framework developed to execute remedial actions. In this section, we proceed on detailing the performed tests to showcase the capabilities and the performance of the presented framework. With this in

mind, this section details the experimental set-up employed to emulate an e-Health scenario as well as the tests performed on top of the testbed.

7.1. E-Health testbed

The experimental testbed is based on an infrastructure deployed within Dell premises in Ovens, Ireland, which replicates the capabilities of the presented e-Health scenario. Figure 8 depicts a schematic of the testbed, highlighting the most important components.

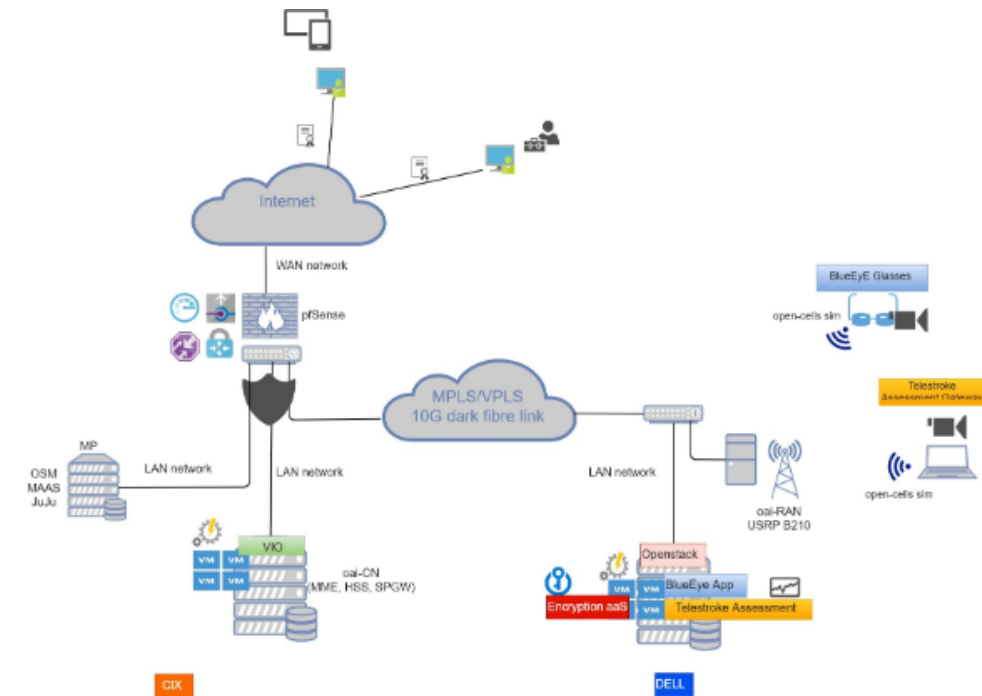


Figure 8: E-Health infrastructure at Dell premises, Ovens, Ireland

On top of the infrastructure, an open source security software, pfSense, is running to provide different security services. Below the firewall and security services, the infrastructure is spanning across 3 racks. A first rack is dedicated to running management services such as OpenSource MANO (OSM) and the software artifacts presented in the proposed architecture; a second rack is running Core network components and cloud operating system, running VIO VMware. A third rack is dedicated to the edge running OpenStack. This edge Virtual Infrastructure Manager (VIM) is hosting the e-Health services. Finally, two stand-alone evolved node B (eNBs) are running

the OpenAirInterface (OAI) software to provide the RAN network segment. To have access to this eNBs, open-cells Subscriber Identity Modules (SIMs) are reprogrammed with uicc/sim programming software provided by open-cells. These SIMs are the ones that emulate the UEs associated to e-Health service, such as wearables inside a roaming ambulance. Given this testbed, the following sub-sections details the specific experiments done to test both the ML model and the Actuation Framework.

7.2. QoP prediction model evaluation

The QoP prediction model aims to train from data that are collected from the vertical service perception. Several real traces are captured having between 2 to 4 hours of length. The total number of instances is 11820. For each instance, seven QoP metrics are considered for determining the quality of the network signal perceived by the ambulance as described in Table 1. Among the real traces, 75% have been used for the training phase and the rest is kept for the testing phase.

As detailed in Section 5, the first step aims to apply a sliding window that will transform the data into a 3D matrix. Each window will observe the KPIs for the last five minutes and it will label the data by observing the future five minutes. The window will advance each 5 seconds. The labeling aims to define the perceived signal strength label for each instance by according thresholds to RSRP, RSRQ and SNR metrics. The result of the previous step is a 3D matrix composed of 11820 instances. Each instance is defined by 7 KPIs where each KPI is defined by 300 values that correspond to 5 minutes of observation. The resulted training set is imbalanced. After applying SMOTE technique, 25 % of all instances correspond to problems in signal strength.

In order to evaluate the model, two techniques are used: (1) Evaluation using cross-validation technique over the training set; and (2) evaluation over the test set. Table 2 presents the results over the training data using a cross validation with 5 folds where the classification is achieved using a random forest. It also illustrates the results over the test data. Table 2 illustrates the numerical results of our experiment where four performance metrics [22] are used: (1) "Accuracy": is the ratio of correctly predicted data over total data; (2) "Recall": is the fraction of the elements of a class that are successfully predicted; (3) "Precision": is the fraction of correct predictions among all the predictions of a class; (4) "F1-score": is a harmonic mean of precision and recall.

Table 2: QoP forecasting model results

Metric	Results with cross-validation	Results over the test set
Accuracy	96%	99%
Precision	91%	99%
Recall	93%	99%
F1-score	93%	99%

As shown in Table 2, the results are promising since it indicates that 93% of the signal strength degradation in the data is detected and 91% of the detected problems are not false alarms. Once the model predicts low perceived signal strength, the Actuation Framework will trigger the event in order to apply a remedial action and to solve the problem before it occurs. This offers the opportunity to the vertical to supervise the performance of its slice and to express its feedback to which SliceNet framework will react accordingly.

7.3. Actuation mechanisms experimentation

The presented QoP prediction model is devoted to predict future perceived degradation in signal quality at the RAN segment in which E2E NS are being supported. Communications over the RAN are of paramount importance, since they enable the necessary mobility of the UE over a geographical area which is a key requirement of e-Health services. As such, proper actuation mechanisms must be proposed and tested in order to react to the predicted degradation.

In the framework of the presented work, in order to overcome potential anomalies in the RAN segment of the NS, two specific actuations have been developed. A first one consist on the handover of the UEs present in the ambulance to another RAN slice that is available in the physical infrastructure. This stems from the fact that signal degradations may be caused due to the mobility of the end-points connected to the e-Health NS, leading to attenuations, delays and other impairments due to the geographical surroundings of the area in which the ambulance is traveling. This being the case, a remedial action entails that the UEs connected to the RAN slice supported over a particular eNB should be re-associated to another RAN slice of another eNB, closer to the current location of the ambulance (i.e. handover).

The second considered actuation consists on enforcing changes on the provisioned resources of the RAN slice, more in concrete, the allocated band-

width. Bandwidth in the RAN is provisioned by reserving a set of contiguous Radio Resource Blocks (RRBs), which provide for the enough network capacity to support the communications over the configured NS. In this regard, it may happen that due to the dynamic nature of the data streams (audio, video) produced by the UEs at the ambulance, in certain points of time the provisioned RRBs are not enough to withstand these data streams, leading to poor video/audio quality and thus affecting the QoE of users. This also translates to some of the QoP metrics being affected (the ones related to the bit-rate), that the QoP prediction model is able to detect. To overcome bandwidth limits, the Actuation Framework engages with the rest of the control/management elements in order to increase the provisioned RRBs at the RAN slices so as to maintain optimal QoE/QoP levels. In both cases of actuation, a policy is placed that ties the output of the ML model to the desired action, handover or bandwidth increase. The setting of the proper policy can be made by a system administrator or by some automated method.

For the developed work in this paper, the FlexRAN SDN controller, which is part of the Mosaic 5G project, has been employed [23]. FlexRAN exposes a REST API that allows external entities to access to the control and configuration operations of an underlying RAN segment. For the presented work, FlexRAN exposes two operations of interest that allow for the materialization of the defined actuations. The first operation of interest allows to change the association of a UE attached to a particular RAN slice in a specific eNB to another RAN slice of the same eNB or other eNB. By exploiting this operation, it is possible to effectively produce a handover of the UEs in the ambulance across RAN slices. The other operation of interest entails the posting of a new configuration for an existing RAN slice, covering aspects such as priority, sharing, scheduling and, more important, the number and position of the RRBs for both up and down links of the slice. Through this operation it is possible to change the RRBs of the slice, thus effectively increasing the bandwidth that is allocated for the communications that originate at the UEs.

In order to provide a Proof of Concept (PoC) of the developed actuation system, we focus on the case of the bandwidth increase actuation. For that, an experimental scenario has been prepared on top of the previously described e-Health testbed. Figure 9 depicts the employed set-up.

The set-up consists on an OSA instance and a P&P control instance that allow collecting the UE-related QoP metrics described previously. A specific Monitoring plugin deployed within the P&P control instance exposes

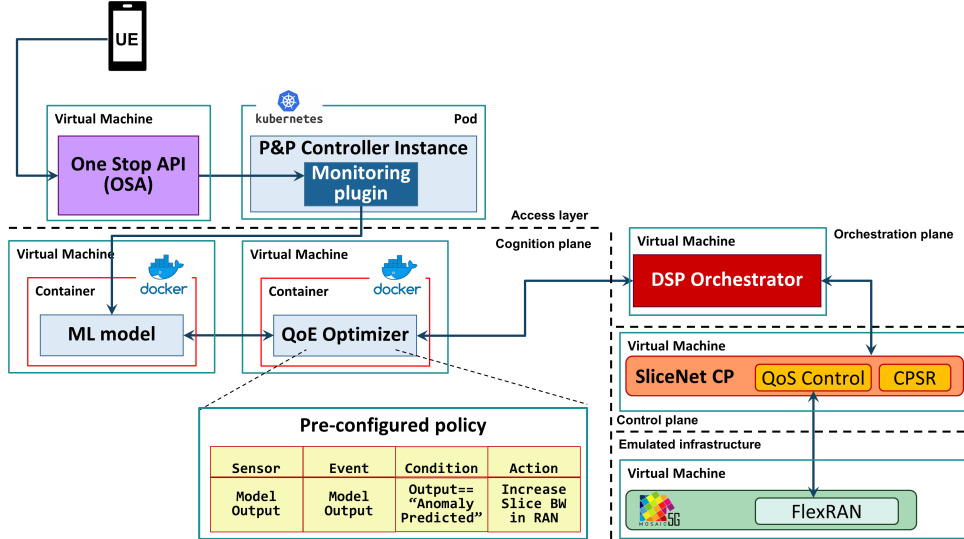


Figure 9: Experimental set-up for actuation PoC.

the monitoring operations towards the UE so as to correctly fetch the desired metrics. Then, these metrics are fed to an instance of the Anomaly Detection ML model which, given the inputs, produces an output that directly feeds and instance of the QoE Optimizer. A pre-configured policy is in place in the QoE Optimizer which dictates that, in the case that the output states that an anomaly is going to happen, it is necessary to increase the bandwidth in the RAN segment of the E2E NS. For the enforcement of the actuation, the QoE Optimizer contacts an instance of the DSP orchestrator. For the particular PoC, the orchestration system has been simplified, with only having one level of orchestration that directly contacts the control functions responsible for the (re-)configuration of slice parameters. The function responsible for that is the QoS Control function within the SliceNet control plane, which allows to enforce any change of QoS parameters (such as bandwidth) of deployed slices by properly engaging with the underlying controls of the physical infrastructure. Lastly, the RAN infrastructure has been emulated employing the aforementioned FlexRAN SDN controller, which, aside of emulating the RAN segment, exposes the control operations necessary to modify the characteristics of RAN slices.

Following the depicted set-up, the PoC consisted on simulating a signal quality degradation to the RAN slice in which the UEs of the ambulance

are connected. To this end, as a first step, a RAN slice has been deployed employing the operations exposed through FlexRAN. The slice has been assigned with a 10% of the RRBs of the total RAN spectrum. Figure 10 depicts in a Java Script Object Notation (JSON) format the details of both the down- and up-link of the slice as exposed through the monitoring API of FlexRAN and captured through wireshark. The highlighted fields are both the percentage of bandwidth occupied by the slice as well as the positioning of the RRBs.

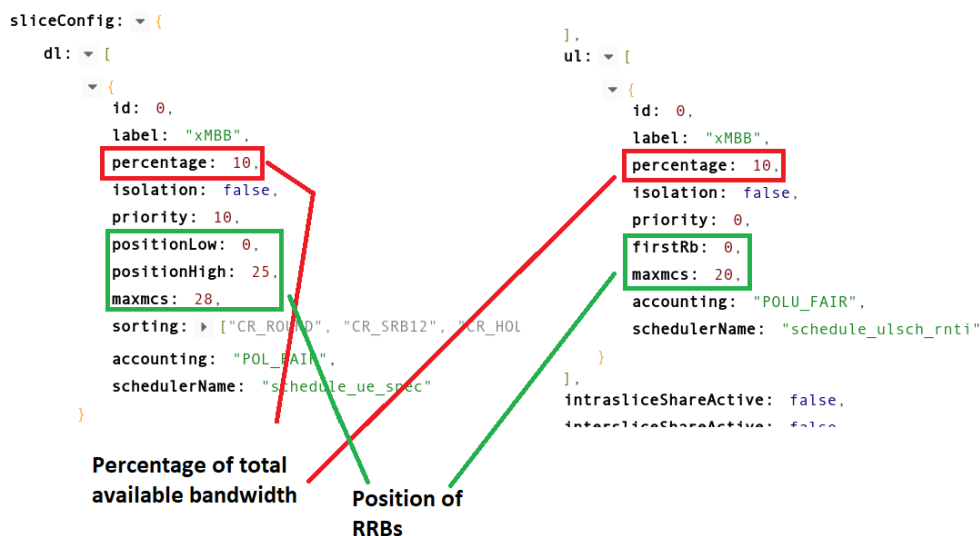


Figure 10: Configured initial RAN slice.

Then, the UE is connected to the deployed slice. Given the configuration of resources and the monitored QoP metrics, the Anomaly Detection ML model detects that a potential anomaly may happen due to the insufficient bit-rate that the UE is experiencing, as a result of the low number of RRBs that have been provisioned. Once this is detected, the full actuation workflow is triggered, which results in a request of increasing the amount of assigned bandwidth of the deployed RAN slice. This is ultimately enforced by the QoS Control module, which contacts the FlexRAN SDN controller specifying the new configuration following the format of the exposed operation. Figure 11 depicts the JSON of the modification operation, stating that the slice needs to be provisioned with 25% of the available RAN spectrum instead of the initial 10%.

```

{
  "dl": [
    {
      "id": 0,
      "percentage": 25
    }
  ],
  "ul": [
    {
      "id": 0,
      "percentage": 25
    }
  ],
  "intrasliceShareActive": true,
  "intersliceShareActive": true
}

```

Figure 11: New RAN slice configuration as a result of the actuation.

Once the FlexRAN controller is contacted, it then takes care of applying the necessary re-configurations at the RAN segment so as the indicated slice is provisioned with the new specified bandwidth. Figure 12 depicts the characteristics of the slice after the re-configuration process. Indeed, it can be appreciated how the desired increase of bandwidth is reflected in its new characteristics by an increased percentage of the provisioned RRBs. After this process, the communications that are supported over the RAN slice find themselves with enough bit-rate to be supported, resulting in the anomaly being adverted.

```

sliceConfig: {
  dl: [
    {
      id: 0,
      label: "xMBB",
      percentage: 25,
      isolation: false,
      priority: 10,
      positionLow: 0,
      positionHigh: 25,
      maxmcs: 28,
      sorting: ["CR_ROI", "CR_SRB12", "CR_HOL"],
      accounting: "POLU_FAIR",
      schedulerName: "schedule_ue_spec"
    }
  ],
  ul: [
    {
      id: 0,
      label: "xMBB",
      percentage: 25,
      isolation: false,
      priority: 0,
      firstRb: 0,
      maxmcs: 20,
      accounting: "POLU_FAIR",
      schedulerName: "schedule_ulsch_rnti"
    }
  ],
  intrasliceShareActive: false,
  intersliceShareActive: false
}

```

Figure 12: Characteristics of the re-configured RAN slice.

8. Conclusion and Future works

Preserving a top communication quality in 5G network slices between the paramedic and the medical teams is a necessity for eHealth verticals. In this regard, a framework has been proposed in this paper allowing the vertical to supervise its slice and to express its feedback regarding the perceived quality of services. Thanks to the developed architecture, it is possible to apply remedial actions to overcome situations in which the quality of the slice, and thus, the supported services, is compromised.

We have shown how the developed ML for QoP estimation is able to digest data coming from UE equipment and transform it into QoP estimations that indicate the quality of the communication over the provisioned E2E slice. From the results, it can be seen that the performance of the model is promising, with high levels of accuracy, showcasing its value as an advanced monitoring function to allow for the quality-aware and cognition-enabled management of e-Health network slices. In addition, we showcased how the developed framework is able to react to outputs coming from the ML model and trigger the necessary remedial actions to overcome the forecast anomalies. We showed how the process is automatically triggered and is able to engage from top to bottom to all involved layers to make effective the necessary changes at the infrastructure level. All in all, the presented experiments show that the proposed architecture is a valuable framework that enables 5G systems with the cognitive management network slices, catering to the special requirements and challenges posed by e-Health vertical services.

As future lines of research, additional prediction models to detect anomalies not only at the RAN segment but also at other segments involved at the E2E slice deployment, such as core, may be investigated. This would allow to have a more precise way of tackling down the network segment responsible for bad quality levels at the service. In addition to that, a method to automate the choice of the corrective action and of the choice of the corresponding policy, not based on static policies disseminated from a policy framework, but rather elaborated based on gathered knowledge and monitoring data, can be investigated. This would allow to adapt the remedial actions to the most updated state of the underlying physical infrastructure and the deployed slices. One way to achieve this automation is by using re-enforcement learning or optimisation algorithms, which could then feed the policy framework for the dynamic creation of policies.

9. Acknowledgment

This work has been funded in part through the European Union's H2020 program, under grant agreement No 761913: project SliceNet. The authors would like to thank all SliceNet partners for their support in this work.

References

- [1] Ericsson Mobility Report, Tech. rep. (june 2019).
URL <https://www.ericsson.com/49d1d9/assets/local/mobility-report/documents/2019/ericsson-mobility-report-june-2019.pdf>
- [2] 3GPPP, 28.801: Telecommunication management; Study on management and orchestration of network slicing for next generation network, Tech. rep. (2017).
- [3] ITU-T, P.10/G.100: Vocabulary and effects of transmission parameters on customer opinion of transmission quality, Tech. rep. (2017).
- [4] ITU-T, E.800: Definitions of terms related to quality of service, Tech. rep. (2009).
- [5] NGMN, 5G White Paper, Tech. rep. (2015).
- [6] P. L. Vo, M. N. H. Nguyen, T. A. Le, N. H. Tran, Slicing the edge: Resource allocation for ran network slicing, *IEEE Wireless Communications Letters* 7 (6) (2018) 970–973.
- [7] S. Husain, A. Kunz, A. Prasad, K. Samdanis, J. Song, Mobile edge computing with network resource slicing for internet-of-things, in: 2018 IEEE 4th World Forum on Internet of Things (WF-IoT), 2018, pp. 1–6.
- [8] SelfNet: A Framework for Self-organized Network Management In Virtualized and Software Defined Networks.
URL <https://selfnet-5g.eu/>
- [9] CogNet – Building an Intelligent System of Insights and Action for 5G Network Management.
URL <https://5g-ppp.eu/cognet/>

- [10] A. Martin, J. Egaña, J. Flórez, J. Montalbán, I. G. Olaizola, M. Quartulli, R. Viola, M. Zorrilla, Network resource allocation system for qos-aware delivery of media services in 5g networks, *IEEE Transactions on Broadcasting* 64 (2) (2018) 561–574.
- [11] W. Kellerer, P. Kalmbach, A. Blenk, A. Basta, M. Reisslein, S. Schmid, Adaptable and data-driven softwarized networks: Review, opportunities, and challenges, *Proceedings of the IEEE* 107 (4) (2019) 711–731.
- [12] V. Sciancalepore, F. Z. Yousaf, X. Costa-Perez, z-torch: An automated nfv orchestration and monitoring solution, *IEEE Transactions on Network and Service Management* 15 (4) (2018) 1292–1306.
- [13] Y. Syu, C. Wang, Y. Fanjiang, Modeling and forecasting of time-aware dynamic qos attributes for cloud services, *IEEE Transactions on Network and Service Management* 16 (1) (2019) 56–71.
- [14] V. Chandrasekhar, Y. Heng, J. Cho, J. Lee, J. Zhang, J. G. Andrews, Experience-centric mobile video scheduling through machine learning, *IEEE Access* 7 (2019) 113017–113030.
- [15] M. Torres Vega, C. Perra, F. De Turck, A. Liotta, A review of predictive quality of experience management in video streaming services, *IEEE Transactions on Broadcasting* 64 (2) (2018) 432–445.
- [16] SliceNet: End-to-End Cognitive Network Slicing and Slice Management Framework in Virtualised Multi-Domain, Multi-Tenant 5G Networks. URL <https://slicenet.eu/>
- [17] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, Smote: Synthetic minority over-sampling technique, *J. Artif. Int. Res.* 16 (1) (2002) 321–357.
- [18] J. O. Ramsay, B. W. Silverman, *Functional data analysis*, 2nd Edition, Springer Series in Statistics, Springer, New York, 2005.
- [19] Y. B. Slimen, S. Allio, J. Jacques, Anomaly prevision in radio access networks using functional data analysis, in: *2017 IEEE Global Communications Conference, GLOBECOM 2017*, Singapore, December 4-8, 2017, IEEE, 2017, pp. 1–6.

- [20] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning, Springer Series in Statistics, Springer New York Inc., New York, NY, USA, 2001.
- [21] IETF, SUPA Policy-based Management Framework, Tech. rep. (2017).
- [22] N. Japkowicz, M. Shah, Evaluating Learning Algorithms: A Classification Perspective, Cambridge University Press, 2011.
- [23] FlexRAN: First Open-source Implementation of a Flexible and Programmable Platform for Software-Defined Radio Access Networks.
URL <http://mosaic-5g.io/flexran/>