

RESEARCH

Open Access

A model for virtual radio resource management in virtual RANs

Sina Khatibi^{1,2*} and Luis M Correia^{1,2}

Abstract

The combination of Network Function Virtualisation (NFV) and cloud-based radio access network (C-RAN) is a candidate approach for the next generation of mobile networks. In this paper, the novel concept of virtual radio resources, which completes the virtual RAN paradigm, is proposed. The key idea is to aggregate (and manage) all the physical radio resources, to create virtual wireless links, and to offer Capacity-as-a-Service. Due to the isolation among instances, network element abstraction, and a multi-radio access techniques (RAT) structure, the virtualisation approach leads to relatively more efficient and flexible RANs than former ones. Virtual network operators (VNOs) ask for wireless connectivity in the form of capacity per service, hence, not dealing with physical radio resources at all. A model for the management of virtual radio resources is proposed, which can even support the shortage of resources. A practical heterogeneous cellular network is considered as a case study, and results are presented, showing how the virtual radio resource management allocates capacity to services of different VNOs, with different service-level agreements (SLAs) and priority when the overall network capacity reduces down to 45% of the initial one.

Keywords: Virtualisation of radio resources; Virtual radio resource management; Radio access networks; Network Function Virtualisation

1 Introduction

Future mobile networks will have to face the rapid growth of mobile data demand [1]. The candidate approach is to use small cell networks with a dense deployment of base stations (BSs); however, traffic varies drastically, both geographically and temporally [2], which creates constraints that are not solved by this approach. The dimensioning of radio access networks (RANs) for busy hours (i.e. the current approach), guarantees the desired performance during that interval, yet it leads to an inefficient resource usage for the remainder of the time, with relatively high capital and operational expenditure (CAPEX and OPEX) costs.

A solution for this matter lays in the ability of adapting RAN during runtime, based on network changes and traffic demands. Hence, flexibility [3] and cost reduction [4] in RANs became the motivation for their implementation in cloud data centres, in order to achieve centralised processing, collaborative radio, real-time cloud computing

[5], and clean RAN systems [6], also known as cloud-based RAN (C-RAN).

Recent studies are focused on the extension of RANs using Network Function Virtualisation (NFV) [7] to add multi-tenancy support, enabling that multiple virtual network operators (VNOs) can be served over the same infrastructure. The concept of a virtualised eNodeB is introduced in [8], by adding an entity, called 'hypervisor', on the top of physical resources, which allocates these resources among various virtual instances. Using the concept of RAN sharing, the air interface resources (i.e. the LTE spectrum) are dynamically divided among various virtual eNodeBs by the hypervisor. In [9], an adaptive allocation of virtual radio resources in heterogeneous networks is analysed, sharing spectrum among VNOs. In [10], the advantage of a virtualised LTE system is investigated by an analytical model for FTP (File Transfer Protocol) transmission. The concept of joint NFV and C-RAN is discussed in [11,12]. This solution, which is called virtual RAN (V-RAN), provides operators with RAN-as-a-Service (RANaaS).

In this paper, the concept of virtualisation of radio resources to achieve virtual wireless links and to have

* Correspondence: sina.khatibi@inov.pt

¹Department of Electrical and Computer Engineering, Instituto Superior Técnico (IST), University of Lisbon, Av. Rovisco Pais, Lisbon 1049-001, Portugal

²INOV-INESC, Rua Alves Redol Lisbon, 1000-029, Portugal

end-to-end virtual networks [13], by aggregating all the physical resources from different radio access techniques (RATs), in order to offer VNOs with a more efficient wireless connectivity is proposed. In this novel methodology, VNOs ask for wireless capacity from a set of physical network providers to serve their subscribers, and they do not have to deal with the physical infrastructure at all. The RAN provider (RANP), owning the physical infrastructure, is capable of offering Capacity-as-a-Service to VNOs. The advantages of RAN virtualisation compared to RAN sharing (where each operator is allocated a portion of spectrum) comes from network element abstraction, isolation among virtual instances, and the ability to support multi-RATs.

A differentiation of these two concepts can be addressed using the analogy presented in [14], where a process on an operating system (OS) is presented as the equivalent of a session in a network. As depicted in Figure 1, the V-RAN and the virtual machine (VM) can be claimed to be a realisation of the corresponding concepts; likewise, RAN sharing is the equivalent of multi-tasking in OSs. In the virtualisation solutions, there is always a virtual manager, such as VMware, offering isolation and abstraction to the upper levels. The offered isolation makes it possible to have multiple instances with different configurations running over the same physical infrastructure, and it relatively reduces the system downtime. The ease of use is the result of the offered abstraction, since virtual instances do not have to deal with physical resources and their complexity.

The novelty of this paper, besides the presented concept of virtualisation of radio resources, is the proposal of an analytical model for virtual radio resource management (VRRM). For a network with multiple RATs, such as GSM, UMTS, and LTE, the model is capable of estimating the overall network capacity based on a given number of the available radio resource units (RRUs) from each RAT. It also shares the available capacity among the different services of the VNOs, the allocation being based on the VNOs' service level agreements

(SLAs), in which VNOs may be guaranteed with a minimum as well as a maximum capacity per service, or simply served in a best-effort approach. The presented VRRM model satisfies SLAs when there are enough RRUUs and minimises SLA violations in resource shortage cases; in both cases, fairness of resource allocation is considered. In addition to the proposition of the novel VRRM model, an architecture for a V-RAN based on a C-RAN infrastructure, with its required modifications to support virtualisation of radio resources, is briefly addressed.

The rest of this paper is organised as follows. Section 2 presents the V-RAN architecture, Section 3 is about the modelling of the problem of management of virtual radio resources. Section 4 describes the details of the scenario, based on which the proposed model is evaluated; numeric results are being discussed in Section 5. Finally, conclusions are presented in the final section of this paper.

2 Virtual radio access network architecture

In this section, the architecture for a V-RAN using virtualisation of radio resources is discussed. It is based on a C-RAN, with modifications to support Capacity-as-a-Service, as depicted in Figure 2. The key architectural elements are as follows:

- VNOs: network operators that do not own a RAN infrastructure. They ask the virtualisation platform for wireless connectivity in terms of capacity to carry various services traffic with various quality-of-service (QoS) requirements to/from their subscribers.
- Backhaul transport network: a low latency optical transport network, which connects the operators' cores to the physical infrastructure of a RAN.
- Virtualisation platform: the key difference between a C-RAN and a V-RAN. On the one hand, it is in charge of abstracting the physical infrastructure for the VNOs, while on the other hand, it handles the request of VNOs through the available physical

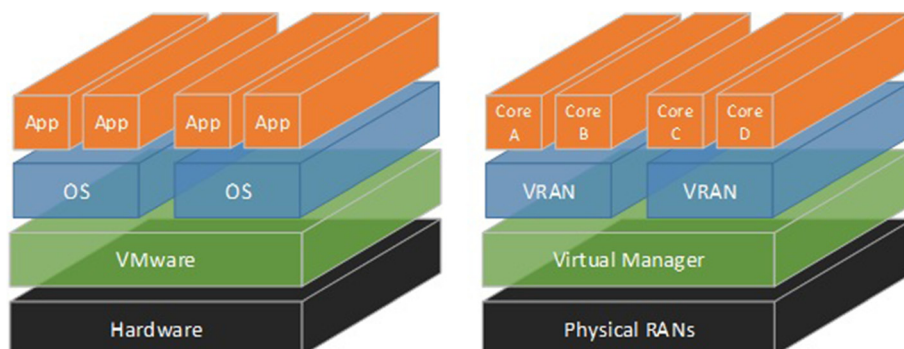


Figure 1 Comparison between V-RAN and VM.

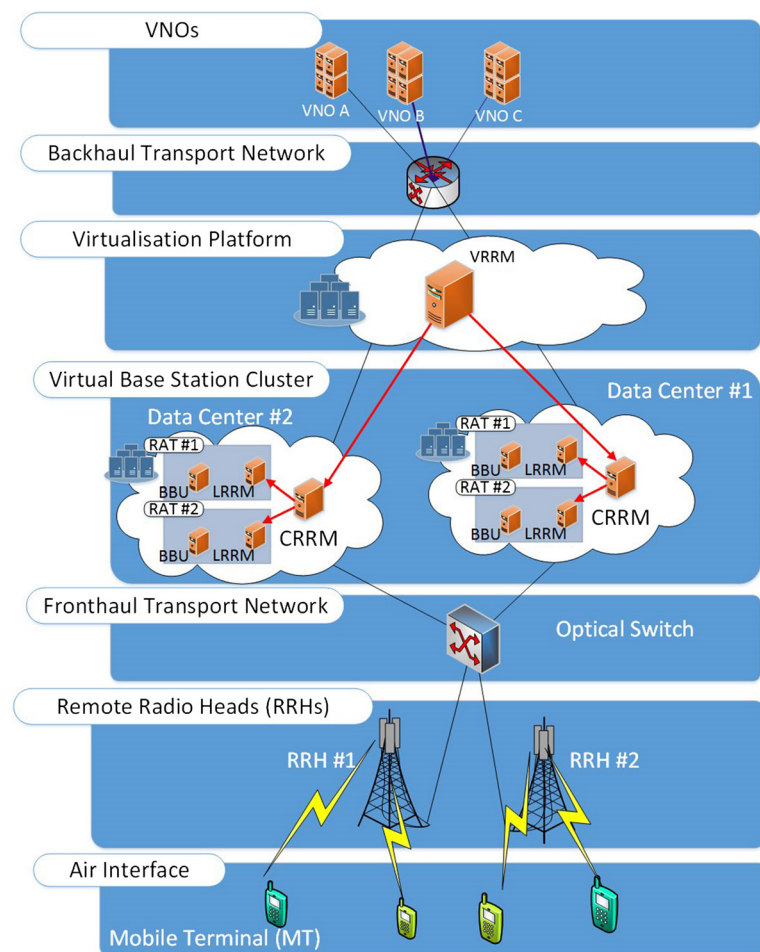


Figure 2 Architecture of V-RAN.

resources. The most important functionality of the virtualisation platform is the VRRM, the highest manager, which is in charge of translating VNO requirements and SLAs through sets of policies onto the lower levels. It optimises the usage of virtual radio resources without dealing with the management of physical resources. Nevertheless, reports and monitoring information (e.g. estimated remained capacity) received from lower levels enable it to improve policies.

- BBU (baseband units) pools' data centre: a set of VMs used for baseband processing of traffic among user terminals and network cores.
- Fronthaul transport network: it transmits digitised radio signals between BBU pools and remote radio heads (RRH), using Common Public Radio Interface (CPRI) with high data rates over optical fibres. The optical equipment needs to have the lowest delay possible, since the maximum round-trip delay must be below 150 μ s (i.e. a maximum of 15 km of BBU-RRH distance) [15]. The optical switch is a

non-smart manageable switch, enabling the scaling or the migration of BBU pools among multiple data centres.

- RRHs: the transceivers in charge of exchanging data and control traffic to/from mobile terminals (MTs) through the air interface, supporting multiple RATs.

By comparing C-RANs with current mobile networks, it can be seen that the eNodeB has been divided into RRHs, fibre optics, and BBU pools. The virtualisation platform, which offers isolation, element abstraction, and multi-tenancy, does not exist in current networks. The changes in the architecture and the dedicated hardware replacement by VMs in data centres provide high flexibility, resource efficiency, and cost reduction.

3 Modelling of virtual radio resource management

The management hierarchy of virtual radio resources is also shown in Figure 2, consisting of VRRM on the top of the usual radio resource management entities in heterogeneous access networks [16], common RRM

(CRRM), and local RRM (LRRM). The VRRM estimates the total network data rate, then, having the available estimated resources, it allocates capacity to the different services of each VNO so that minimum and maximum guaranteed capacities are met. This section presents the analytical modelling, considering the estimation of resources, and their allocation without and with violation of SLAs.

3.1 Estimation of available resources

In general, the data rate of an RRU assigned to an MT varies between zero and the maximum data rate based on various parameters, e.g. RAT, modulation, and coding schemes. Therefore, it can be given as a function of channel quality, i.e. signal to noise ratio (SNR), as follows:

$$R_{b_{\text{RAT}_i}[\text{Mbps}]}(\rho_{\text{in}}) \in [0, R_{b_{\text{RAT}_i}[\text{Mbps}]}^{\max}] \quad (1)$$

where

- $R_{b_{\text{RAT}_i}}$ is the data rate of a single RRU of the i th RAT,
- ρ_{in} is the input SNR, and
- $R_{b_{\text{RAT}_i}}^{\max}$ is the maximum data rate of a single RRU of the i th RAT.

In [17], a heterogeneous cellular network is modelled as a K -tier network, where each tier models the BSs of a particular class. It is assumed that the BSs in a given tier are spatially distributed as a Poisson point process (PPP) with a given density and transmission power. The received power is assumed to be exponentially distributed (i.e. Rayleigh fading is assumed for the signal magnitude). It is shown that the cumulative distributed function (CDF) of the input SNR for an interference limited network, where MTs are connected to the BS with the strongest signal, can be written as follows:

$$P_{\rho_{\text{in}}[\text{dB}]}(\rho_{\text{in}}) = 1 - e^{-\frac{0.2}{\alpha_p} \ln(10) \rho_{\text{in}}[\text{dB}]} \quad (2)$$

where

- α_p is the path loss exponent (which values are larger or equal to 2).

Based on real logs, the data rates of different access technologies, as a function of input SINR and vice versa, have been presented in [18]. In the next step, for the sake of simplicity, these functions have been approximated by an equivalent polynomial of degree 5; hence,

the SINR can be written as a function of data rate as follows:

$$\rho_{\text{in}}[\text{dB}](R_{b_{\text{RAT}_i}}) = \sum_{k=0}^5 a_k \left[\frac{\text{dB}}{\text{Mbps}^k} \right] R_{b_{\text{RAT}_i}[\text{Mbps}]}^k \quad (3)$$

where

- a_k are coefficients of a polynomial approximation of SINR as a function of data rate in each access technology listed in [13].

By substituting this polynomial in (2) and adding the boundary conditions addressed in (1), the CDF of a single RRU of RAT_i is:

$$P_{\text{Rb}}(R_{b_{\text{RAT}_i}[\text{Mbps}]}) = \frac{e^{-\frac{0.2}{\alpha_p} \ln(10) a_0} - e^{-\frac{0.2}{\alpha_p} \ln(10) \sum_{k=0}^5 a_k (R_{b_{\text{RAT}_i}})^k}}{e^{-\frac{0.2}{\alpha_p} \ln(10) a_0} - e^{-\frac{0.2}{\alpha_p} \ln(10) \sum_{k=0}^5 a_k (R_{b_{\text{RAT}_i}}^{\max})^k}} \quad (4)$$

In the next step, the overall capacity of a RAT is estimated as follows:

$$R_{b_{\text{tot}}}^{\text{RAT}_i} = \sum_{n=1}^{N_{\text{RRU}}^{\text{RAT}_i}} R_{b_n}^{\text{RAT}_i} \quad (5)$$

where

- $N_{\text{RRU}}^{\text{RAT}_i}$ is the number of RRUs in the i th RAT,
- $R_{b_{\text{tot}}}^{\text{RAT}_i}$ is the data rate from a i th RAT pool, and
- $R_{b_n}^{\text{RAT}_i}$ is the data rate from the n th RRU of the i th RAT.

Based on [19], the probability density function (PDF) of each RAT, assuming the RRUs are independent, is equal to the convolution of all the PDFs of that RAT's RRU. From (4), the PDF of a single RRU is calculated (and then numerically sampled with a step of 10 kbps). To compute the total data rate PDF of each RAT, the PDF of the entire RRUs is convolved.

The resource pools of RATs can be aggregated under the supervision of the CRRM, and the total data rate from all RATs is the summation of the total data rate from each of them:

$$R_{b[\text{Mbps}]}^{\text{CRRM}} = \sum_{i=1}^{N_{\text{RAT}}} R_{b_{\text{tot}}[\text{Mbps}]}^{\text{RAT}_i} \quad (6)$$

The PDF of the total network data rate is computed by convolving all the RATs' PDFs. By obtaining the total network CDF and PDF, an estimation of available

network capacity is in hand to be used in the allocation procedure, as described in the next subsections.

3.2 Allocation of resources

After estimating the total network capacity, the VRRM has to allocate it to the various services of the networks. The key objective in the allocation of resources is to maximise the usage efficiency in addition to meeting the constraint set. The algorithms of resource allocation have also to consider the priority of the different services of different VNOs based on their SLAs. For instance, conversation (e.g. VoIP) and streaming (e.g. video streaming) service classes are delay sensitive, but they have almost constant data rates. The allocation to these services of data rates higher than contracted capacities do not increase quality of service (QoS), in contrast to interactive (e.g. FTP) and background (e.g. email) service classes, where the increase of data rates can indeed improve a user's quality of experience (QoE); hence, operators offering the former services are not interested in allocating higher data rates. Based on the service set and requirements, VNOs may have different SLAs, but these SLAs can generally be categorised into three types of contract:

- **Guaranteed bit rate (GB)**, in which the VNO is guaranteed a minimum as well as a maximum level of data rates, regardless of the network status. In other words, the total satisfaction of the VNO is achieved when the maximum guaranteed data rate is allocated to it. The upper boundary in this type of SLA enables VNOs to have full control on their networks. For instance, a VNO offering VoIP to its subscribers may foresee to offer this service to only 30% up 50% of its subscribers simultaneously. The VNO can put this policy into practice by choosing a guaranteed SLA for its VoIP service. It is expected that subscribers always experience a good QoS in return of relatively more expensive services.
- **Best effort with minimum guaranteed (BG)**, where the VNO is guaranteed with a minimum level of service, but the request for higher data rates than the guaranteed one is served in a best-effort manner. In this case, although VNOs do not invest as much as former ones, they can still guarantee the minimum QoS to their subscribers. From the subscribers' viewpoint, the acceptable service (not as good as the other ones) is offered with a relatively lower cost.
- **Best effort (BE)**, in which the VNO is served in a pure best-effort manner. Operators, and consequently their subscribers, in return, may suffer from low QoS and resource starvation during busy hours.

Hence, the objective function is formed as the total weighted network data rate:

$$f_{\mathbf{R}_b}^{\text{cell}}(\mathbf{R}_b) = \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} W_{ji}^{\text{Srv}} R_{bji}^{\text{Srv}} \quad (7)$$

where

- R_{bji}^{Srv} is the serving data rate for service j of VNO i ,
- \mathbf{R}_b is the vector of serving data rates,
- N_{VNO} is the number of served VNOs,
- N_{srv} is the number of services for each VNO, and
- W_{ji}^{Srv} is the weight of serving unit of data rate for the j th service of the i th VNO, where $W_{ji}^{\text{Srv}} \in [0, 1]$.

The weights in (7) are used to prioritise the allocation of data rates to different services of different VNOs. Services with the higher weights are served with the higher data rates. The choice of these weights is based on the SLAs between the VNOs and VRRM.

There are also constraints in the allocation of data rates that should not be violated. The fundamental constraint is the total network capacity, estimated in the previous subsection. The summation of the entire assigned data rate to all services should not be greater than the total estimated capacity of the network:

$$\sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} R_{bji}^{\text{Srv}} \leq R_b^{\text{CRRM}} \quad (8)$$

where

- R_b^{CRRM} is the estimated total data rate that can be provided by CRRM from various RATs.

However, the optimisation of this objective function in the current situation may not lead to a desirable situation: services with the highest serving weight receive almost all the resources, while the other services are allocated by the minimum possible data rate; this way of resource allocation is neither fair nor desirable. In contrast, the ideal case is when the normalised data rate (i.e. the data rate divided by the serving weight) of all services, and consequently the normalised average, has the same value. This can be expressed as follows:

$$\frac{R_{bji}^{\text{Srv}}}{W_{ji}^{\text{Srv}}} - \frac{1}{N_{\text{VNO}} N_{\text{srv}}} \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} \frac{R_{bji}^{\text{Srv}}}{W_{ji}^{\text{Usg}}} = 0 \quad (9)$$

Nevertheless, resource efficiency and fair allocation are two contradict goals. For instance, if one assumes a network with a 100-Mbps capacity to serve two services with serving weights of 0.8 and 0.2, by considering only (7), all the network capacity has to be allocated to the first service

(the one with a serving weight of 0.8), while a fair allocation is achieved when the first service receives 80 Mbps and the other 20 Mbps. As expected, the increment of the data rate in one of them leads to the decrement in the other; hence, instead of having the fairest allocation possible (i.e. the deviation of all normalised data rates from the normalised average is zero), the minimisation of the total deviation from the normalised average is used:

$$f_{\mathbf{R}_b}^{fr} = \min_{\mathbf{R}_{b_{ji}}^{Srv}} \left\{ \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} R_{b_{ji}}^D [\text{Mbps}] \right\} \quad (10)$$

where

- $f_{\mathbf{R}_b}^{fr}$ is the fairness objective function and
- $R_{b_{ji}}^D$ is the deviation from the normalised average for service j of VNO i , given by the following:

$$R_{b_{ji}}^D [\text{Mbps}] = \left| \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} - \frac{1}{N_{VNO} N_{srv}} \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} \right| \quad (11)$$

In order to convert the problem into a linear form, (11) can be written as follows:

$$\begin{aligned} -R_{b_{ji}}^f [\text{Mbps}] &\leq \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} - \frac{1}{N_{VNO} N_{srv}} \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} \\ &\leq R_{b_{ji}}^f [\text{Mbps}] \end{aligned} \quad (12)$$

where

- $R_{b_{ji}}^f$ is the boundary for deviation data rate from the normalised average for service j of VNO i .

According to (12), $R_{b_{ji}}^f$ is always larger or equal to $R_{b_{ji}}^D$, and minimising the former leads to the minimisation of the latter. Therefore, (10) reformulated into a form originated from [20] as follows:

$$\begin{aligned} f_{\mathbf{R}_b}^{fr} &= \min_{\mathbf{R}_{b_{ji}}^{Srv}, \mathbf{R}_{b_{ji}}^f} \left\{ \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} R_{b_{ji}}^f [\text{Mbps}] \right\} \\ \text{s.t.} \quad &\begin{cases} \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} - \frac{1}{N_{VNO} N_{srv}} \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} \leq R_{b_{ji}}^f [\text{Mbps}] \\ -\frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} + \frac{1}{N_{VNO} N_{srv}} \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} \frac{R_{b_{ji}}^{Srv} [\text{Mbps}]}{W_{ji}^{Srv}} \leq R_{b_{ji}}^f [\text{Mbps}] \end{cases} \end{aligned} \quad (13)$$

It is worthwhile noting that the fairness for services with minimum guaranteed data rates applies only to the

amount exceeded over the minimum guaranteed level. As the network capacity increases, the summation in (7) increases as well; therefore, in order to combine it with (13), the fairness intermediate variable, R_{ji}^f , has to adapt to the network's capacity:

$$f_{\mathbf{R}_b}^v(\mathbf{R}_b) = f_{\mathbf{R}_b}^{\text{cell}}(R_{b_{ji}}^{\text{cell}}) - W^f f_{\mathbf{R}_b}^f(\mathbf{R}_b^f) \quad (14)$$

where

- $W^f \in [0, 1]$ is the fairness weight in the objective function, indicating how much weight should be put on the fair allocation and
- \mathbf{R}_b^f is the vector of intermediate fairness variable:

$$\mathbf{R}_b^f = \left\{ R_{b_{ji}}^f | j = 1, 2, \dots, N_{srv} \text{ and } i = 1, 2, \dots, N_{VNO} \right\} \quad (15)$$

- $f_{\mathbf{R}_b}^f$ is the fairness function:

$$f_{\mathbf{R}_b}^f(\mathbf{R}_b^f) = \sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} \left(\frac{R_{b_{ji}}^{\text{CRRM}} [\text{Mbps}]}{R_{b_{ji}}^{\text{min}} [\text{Mbps}]} R_{b_{ji}}^f [\text{Mbps}] \right) \quad (16)$$

where

- R_b^{min} is the minimum average data rate among all the network services (i.e. VoIP).

The division of the network capacity by the minimum average data rate of services gives the maximum possible number of users in the network with a given network capacity and service set. By multiplying the fairness variable by the maximum number of users, the balance of these two objectives (i.e. network throughput and fairness) can be kept.

In addition, there are more constraints for VRRM to allocate data rates to various services, which should not be violated. The very fundamental constraint is the total network capacity estimated in the previous section. The summation of the entire assigned data rates to all services should always be smaller than the total estimated capacity of the network:

$$\sum_{i=1}^{N_{VNO}} \sum_{j=1}^{N_{srv}} R_{b_{ji}}^{Srv} [\text{Mbps}] \leq R_b^{\text{CRRM}} [\text{Mbps}] \quad (17)$$

The offered data rate to the guaranteed and the best effort with minimum guaranteed services imposes the next constraints. The allocated data rates related to these services have to be higher than the minimum guaranteed level (for guaranteed and best effort with minimum

guaranteed) and lower than the maximum guaranteed one (for guaranteed services only):

$$R_{b_{ji}[\text{Mbps}]}^{\text{Min}} \leq R_{b_{ji}[\text{Mbps}]}^{\text{Srv}} \leq R_{b_{ji}[\text{Mbps}]}^{\text{Max}} \quad (18)$$

where

- $R_{b_{ji}}^{\text{Min}}$ is the minimum guaranteed data rate of service j of VNO i and
- $R_{b_{ji}}^{\text{Max}}$ is the maximum guaranteed data rate of service j of VNO i .

3.3 Resource allocation with violation

However, in the allocation process, there are situations where the resources are not enough to meet all guaranteed capacity and the allocation optimisation is no longer feasible. Data centre migration is a practical example of this case. A simple approach in these cases is to relax the constraints by the introduction of violation (also known as slack) variables. In case of VRRM, the relaxed constraint is as follows:

$$R_{b_{ji}[\text{Mbps}]}^{\text{Min}} \leq R_{b_{ji}[\text{Mbps}]}^{\text{Srv}} + \Delta R_{b_{ji}[\text{Mbps}]}^v \quad (19)$$

$$\Delta R_{b_{ji}[\text{Mbps}]}^v \geq 0$$

where

- $\Delta R_{b_{ji}}^v$ is the violation variable for the minimum guaranteed data rate of service j of VNO i .

By introducing the violation parameter, the former infeasible optimisation problem turns into a feasible one. The optimal solution maximises the objective function and minimises the weighted average constraint violations. The weighted average constraint violation is defined as follows:

$$\Delta \bar{R}_b^v[\text{Mbps}] = \frac{1}{N_{\text{VNO}} N_{\text{srv}}} \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v \quad (20)$$

where

- $\Delta \bar{R}_b^v$ is the average constraint violation and
- W_{ji}^v is the weight of violating minimum guaranteed data rate of service j of VNO i , where $W_{ji}^v \in [0, 1]$.

The objective function presented in (14) has also to be changed. The new objective function, the relaxed one, has to contain the minimisation of violations in addition to the maximisation of former objectives. Although the average constraint violation has a direct relation with the allocated data rate to services, where the increment in one leads to the decrement of the other, it does not have the same relation with fairness. It can be claimed

that the maximisation of fairness and minimisation of constraint violations are independent. Therefore, the final objective function considering both issues has to consider the same approach for minimisation of the violations as well as fairness. In other words, the fairness variable is weighted as it is presented in (17) to compensate the summation of weighted data rate of various services. The derivation from fair allocation, which is desired to be as minimum as possible, leads to a relatively higher weight in the objective function and may confiscate the constraint violation strategies. Therefore, the average constraint violation also has to be placed in the objective function in a similar way:

$$f_{\mathbf{R}_b}^v(\mathbf{R}_b) = f_{\mathbf{R}_b}^{\text{cell}}(\mathbf{R}_b^{\text{cell}}) - f_{R_b^v}^{\text{vi}}(\Delta \bar{R}_b^v) - W^f f_{R_b}^f(\mathbf{R}_b^f) \quad (21)$$

where $f_{R_b^v}^{\text{vi}}$ is the constraint violation function:

$$f_{R_b^v}^{\text{vi}}(\Delta \bar{R}_b^v) = \frac{R_{b[\text{Mbps}]}^{\text{CRRM}}}{R_{b[\text{Mbps}]}^{\text{min}}} \Delta \bar{R}_b^v[\text{Mbps}] \quad (22)$$

However, the definition of fairness in a congestion situation is not the same. The fairness objective in the normal case is to have the same normalised data rate for all services. As a reminder, when the network faces congestion, there are not enough resources to serve all services with the minimum acceptable data rates. Therefore, some of best-effort services are not allocated any capacity at all, and some violation is also introduced in the guaranteed data rates. In this case, fairness is to make sure that the weighted violation of all services is the same. The ideal fairness with this approach is as follows:

$$W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v - \frac{1}{N_{\text{VNO}} N_{\text{srv}}} \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v = 0 \quad (23)$$

The violation data rates for the best-effort services are always zero; consequently, (13) is changed to the following:

$$f_{\mathbf{R}_b}^f = \min_{R_{b_{ji}}^{\text{Srv}}, R_{b_{ji}}^f} \left\{ \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} R_{b_{ji}[\text{Mbps}]}^f \right\}$$

$$s.t. \begin{cases} W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v - \frac{1}{N_{\text{VNO}} N_{\text{srv}}} \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v \leq R_{b_{ji}[\text{Mbps}]}^f \\ -W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v + \frac{1}{N_{\text{VNO}} N_{\text{srv}}} \sum_{i=1}^{N_{\text{VNO}}} \sum_{j=1}^{N_{\text{srv}}} W_{ji}^v \Delta R_{b_{ji}[\text{Mbps}]}^v \leq R_{b_{ji}[\text{Mbps}]}^f \end{cases} \quad (24)$$

The management of virtual radio resources is a complex optimisation problem since the network status and constraints vary in time. Among various possible techniques and approaches for solving this problem, partial

VRRM seems to be the simplest one. In this approach, the main optimisation problem is broken into multiple sub-problems. In other words, the time axis is divided into decision windows, and VRRM maximises the objective function in each of these intervals, independently. However, it is worth noting that decisions in each interval affect directly the network state, and the outcome at a certain interval depends on the decisions and states in previous intervals; the optimal solution has to take this dependency into consideration. As a consequence, the output of partial VRRM may only be a local minimum and not the global one. Nevertheless, partial VRRM is a simple solution, which can be used as the starting step and reference point.

Figure 3 illustrates a decision window of VRRM, CRRM, and LRRMs. The VRRM decision window contains multiple CRRM ones, during which CRRM applies the decided policy set. In the next decision window of VRRM, after multiple network stages, the VRRM updates the network situation and makes the new decision for the next time interval.

The aforementioned optimisation problem is solved by MATLAB Linear Programming (LP) problem solver (i.e. linprog function) [21]. The method used in this function is the interior-point LP [22], which is a variant of Mehrotra's predictor-corrector algorithm [23], a primal-dual interior-point method. The termination tolerance on the function is chosen to be 10^{-8} .

4 Scenario

A number of scenarios are chosen to evaluate the performance of the proposed model. The key parameters of these scenarios are cell layout, the RATs' configuration, the VNOs, and the service set.

The RRHs are capable of supporting multiple RATs, which are OFDMA (based on LTE-Advance), CDMA (based on UMTS), and FDMA/TDMA (based on GSM), and their flexibility enables various cell layout for these RATs. The considered layout, illustrated in Figure 4, offers full coverage using TDMA cells with the radius of 1.6 km, CDMA cells with 1.2 km, and OFDMA cells with 0.4 km. It is assumed that the coverage area is divided into serving areas, over which a VRRM is

operating. Dividing the coverage area to different serving areas makes it possible to consider different policies for different regions (e.g. for residential or commercial regions). In these scenarios, the serving area for each VRRM is considered to be as big as the TDMA cell; hence, each serving area is covered by 1 TDMA cell, approximately 1.7 CDMA cells, and 16 OFDMA ones.

The details of each RAT configuration, such as the number of cells and the number of RRU per RAT, are presented in Table 1. For the CDMA cells, in which the serving area covers an area equivalent to area of 1.7 cells, it is assumed that the radio resources are distributed uniformly and the available resources for this RAT are 1.7 times the resources of a single cell. Moreover, variations of the reference scenario are considered, in which the serving area is covered with a lower number of OFDMA cells temporarily. A lower cell number leads to a lower network capacity; hence, network capacity and VRRM performance are compared in these scenarios. The minimum number of OFDMA cells is chosen to be 5, an extreme case where the network capacity is reduced to 45% of the reference scenario's capacity.

Furthermore, 3 VNOs, each one with 300 subscribers, are assumed to operate in this area, and the average required data rate for each of them is 6.375 Mbps [24]. Hence, the contracted data rate for each of these operators is 1,912.5 Mbps. It is worth noting that the choice of the average data rate is just used to consider realistic boundaries for the guaranteed data rates. Although they have the same number of subscribers and contracted data rate, they are different SLAs as follows:

- VNO GB, the allocated data rates for services are guaranteed to be in a range [50, 100]% of the service data rate.
- VNO BG has best effort with a minimum of 25% of service data rate guaranteed SLA.
- Services of VNO BE are served all in a best-effort manner.

All of these VNOs offer the same set of services to their subscribers. These services and their volume share

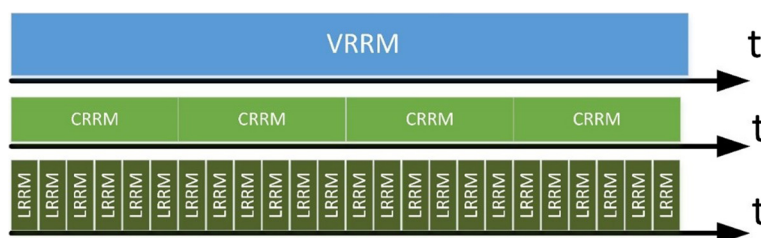
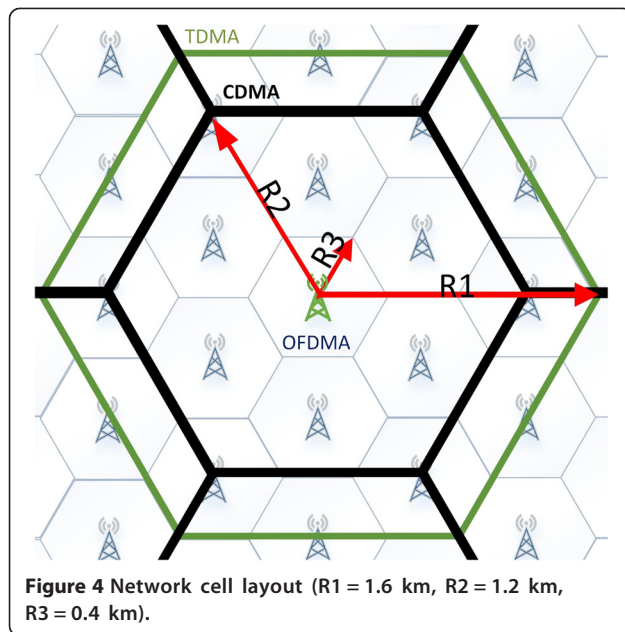


Figure 3 Decision window of VRRM and CRRM.



of an operator traffic are listed in Table 2 which are adopted from [25,26].

Finally, the serving and the violation weights of the services are based on general service classes: conversational (0.4), streaming (0.3), interactive (0.2), and best effort (0.05); in order not to compromise the objective function for having a higher fairness, the fairness weight, W^f , is heuristically chosen to be equal to the lowest serving weight (0.05).

5 Analysis of results

Results for the reference scenario and its variation were obtained, being presented and analysed from three main perspectives: the total network capacity and the capacity of VNOs, the allocated data rate to each service of a VNO in the reference scenario, and finally the allocated data rate to each service class in VNO GB.

5.1 Total network and VNO capacity

The total network capacity of network is achieved by obtaining the PDF of different RATs, as presented in (4). One compares the concept of virtualisation of radio resources and RAN sharing by considering the CDFs of the total network. Since all three VNOs have the same traffic demand, RRUs are divided into three equal parts in RAN sharing, whereas in the V-RAN approach all

Table 1 Different RAT cell radius

RAT	Number cells	System	N_{RRU}^{RAT}	Total RRUs
OFDMA	16	LTE	500	8,000
CDMA	1.7	UMTS	45	80
TDMA	1	GSM	75	75

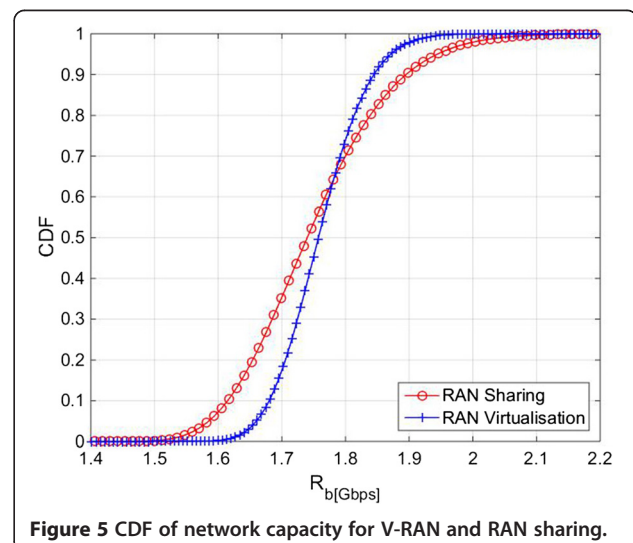
Table 2 Network traffic mixture

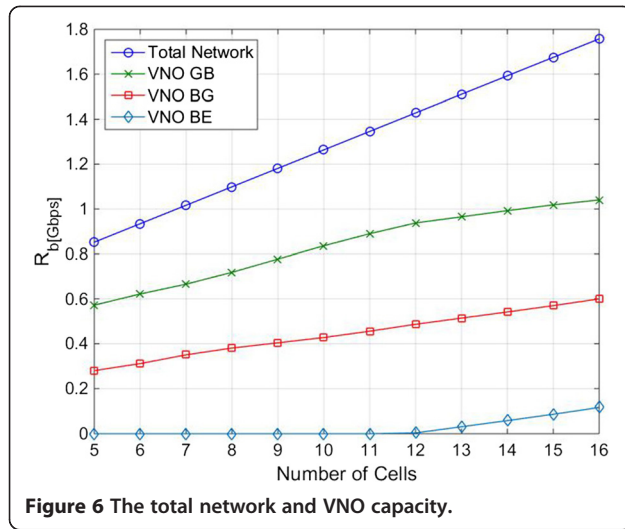
Service		Volume [%]	W_{ji}^{Srv}	W_{ji}^v
VoIP (Vol)		1.0	0.04	0.36
Music streaming (MuS)		3.0	0.03	0.27
File sharing (FTP)		3.5	0.02	0.18
Web browsing (WWW)		11.9	0.02	0.18
Social networking (SoN)		14.4	0.02	0.18
Email (Ema)		1.0	0.005	0.045
M2M	Smart metres (MMM)	1.475	0.005	0.045
	e-Health (MME)	1.475	0.02	0.18
	ITS (MMI)	1.475	0.04	0.36
	Surveillance (MMS)	1.475	0.03	0.27
Mobile video	Video calling (ViC)	2.75	0.04	0.36
	Video streaming (VoS)	56.95	0.03	0.27

RRUs are aggregated. Taking the data from Table 1, the CDF of the total network for RAN sharing and the V-RAN approach using (6) is obtained (Figure 5).

For the sake of simplicity, the RAN sharing CDF using one third of resources is multiplied by three. However, it should be reminded that the total capacity of the network using all the aggregation can be achieved by convolving the PDF of each spectrum slice and not simply summing them. It can be seen that for 50% of the time, the total V-RAN network capacity is 1,800 Mbps, where RAN sharing offers less than 1,782 Mbps. The highest difference can be seen where the CDF is equal to 0.1, in which case, the relative data rate for the V-RAN is 1 725 Mbps, while RAN sharing offers only 1,656 Mbps.

Figure 6 illustrates the total network capacity when different number of cells is used to cover the serving area. The total network capacity with the 16 cells (i.e. the reference scenario) is 1,800 Mbps. It reduces to 48.5% of its initial value (i.e. 872.4 Mbps) when the full





coverage is obtained by only five OFDMA BSs. According to the scenario definition, the total guaranteed data rate is 1,434.37 Mbps, which means that there is enough capacity to serve the guaranteed data rate plus the best-effort services. The layout with 12 cells is the marginal point where the network capacity and the total guaranteed data rate are almost equal; the use of only five cells provides a very low capacity.

Considering the allocated data rate to VNOs, as expected, all the capacity allocated to the VNOs decreases by reducing the number of cells. Capacity reduction has a higher impact on VNO BE (the best-effort operator) comparing to VNO's GE and BG, since the network tries to meet these latter VNOs' guaranteed capacity before serving the best-effort one. When there are 12 cells, VNO BE gets almost no data rate, but the other two VNOs still have a relatively acceptable data rate. In this situation, the total network capacity is still higher than the total guaranteed capacity.

The network capacity shrinks to 1,378.15 Mbps when another cell is reduced, i.e. 11 cells, which is lower than the guaranteed data rate. The violation is inevitable for the cell layout with less than 12 cells. In these situations, the main objective function becomes infeasible and VRRM switches to the objective function with violation, presented in (21). While no capacity is allocated to VNO BE, the other two VNOs share the violation between them. Since the model tries to minimise the weighted average violation, it can be seen that VNO GE always receive a relatively larger portion of the network capacity, since it has a higher guarantee rate.

5.2 Data rate allocation in service level

At the service level, Table 3 presents the allocation of data rates to the services of all three VNOs for the

Table 3 Allocated data rate to services when all the cells are available

Services	$R_{b,i}^{Srv}$ [Mbps]		
	VNO GE	VNO BG	VNO BE
VoIP	19.12	21.47	16.69
Music streaming	40.30	25.96	11.62
File sharing	41.21	24.48	7.74
Web browsing	121.54	64.64	7.74
Social networking	145.44	93.33	7.74
Email	11.50	6.72	1.94
M2M-SM	15.08	8.51	1.94
M2M-eH	20.89	14.32	7.74
M2M-ITS	26.30	23.26	16.69
M2M-SV	24.77	18.19	11.62
Video streaming	556.24	283.93	11.62
Video call	42.83	29.82	16.69

reference case with 16 cells; in these conditions, the VRRM is able to allocate the capacity to all services without violating any constraints. As expected, the highest data rate is allocated to video streaming of VNO GB, since it has the highest guaranteed data rate. The lowest data rates are given to Email and M2M Smart Meter services, since they are background ones with the lowest serving weight. The best demonstration of prioritising the services based on their serving weights can be seen in VNO BE, where there is no minimum guaranteed data rate for the services. The highest capacities belong to VoIP, M2M-ITS, and video calls, which are services from the conversational class with the highest data rates; since these services have the same serving weight, they receive the same capacity. Music, M2M-SV, and video streaming are in the second group, i.e. streaming. Services of the interactive class, i.e. FTP, web browsing, social networking, and M2M-eH, received all 7.74 Mbps. The effect of fairness is very well demonstrated in services of VNO BE; although the services have different serving weights, they are served relatively well based on their serving weight. In addition, services with the same serving weights are allocated the same capacity. For the other two VNOs, the services have different guaranteed capacities, and the fairness effect is not as obvious as in VNO BE.

It is worth noting that Table 3 is also showing an interesting difference of best effort with minimum guaranteed services and guaranteed ones. Guaranteed services are bounded by the maximum capacity, and the allocated capacity cannot go higher than this boundary, while best effort with minimum guaranteed service does not have this limitation. Considering VoIP in VNO GB and VNO BG, it can be seen that the latter is allocated with a higher capacity since this service of VNO GB is

assigned with the maximum capacity; for VNO BE, VoIP shows the effect of fairness among services.

5.3 Allocated data rate to each service class in VNO GB

Finally, to study how different services are affected by the changes of the total network capacity from the VNO GB, four service classes are considered, with different serving and violation weights. Figure 7 illustrates the percentage of violation for each of these services by shutting down cells.

Obviously, there is no violation of guaranteed capacity as long as there are more than 12 cells. However, VRRM has no other choice than start violating the level of data rate guaranteed when there are fewer cells. As a matter of fact, the violations have to start by the service with the lowest violation weight. According to the weights presented in Table 2, background services are the first candidate for violation, since they have the lowest serving and violation weights. When there are only 11 cells, the background traffic violation reaches 100%, which means that it is not allocated any capacity at all. Since these services (i.e. Email and M2M-SV) are low-volume services, even the total violation of their capacity cannot cover the shortage of network capacity; therefore, interactive services, the ones with the second lowest violation weights, have also to be subject of capacity violations. Finally, when there are only five cells (the worst network situation considered in this paper), background and interactive services have to be totally shut down (i.e. 100% violation), video streaming suffering a violating up to 8%. Nevertheless, conversational services, the ones with the highest weights, are served without any violations.

In conclusion, it can be seen that by aggregating all the radio resources, their efficiency use increases. In

addition, the effect of having different SLAs (i.e. best effort or guaranteed) and priority (i.e. serving weights) is demonstrated by means of these results.

6 Conclusions

This paper presents the concept of virtualisation of radio resources as the final step towards an end-to-end virtual network by realising a virtual wireless link. It is suggested to aggregate all the physical radio resources and to have a central management to offer VNO's Capacity-as-a-Service. In this solution, VNOs no longer have to deal with the management of physical infrastructure, but rather to ask for capacity in order to serve their subscribers. Using the proposed technique, not only heterogeneous access networks are shared among multiple VNOs but also the ease of use and VNO specific configuration are achieved as the result of network element abstraction and isolation.

In addition, a model for the management of virtual radio resources with the capability of supporting various situations is proposed. The model takes a number of available RRUUs in different RATs as the input and maps them onto the total network capacity. Having an estimation of the available capacity, the model formulates the allocation problem into a linear optimisation problem. The objective function in this problem is to maximise the weighted throughput of network and fairness among the services of the VNOs. The suggested model also tries to meet all guaranteed service levels while offering fairness. When the model fails to find any feasible solution to serve all guaranteed data rates, due to the shortage of resources, it introduces violations to guaranteed data rates. This approach changes the former infeasible solution to a feasible one, and the model aims at minimising the summation of weighted violations. This way, the services with less importance are facing the violation of guaranteed data rates while the more important ones are served properly.

The proposed model is evaluated in a practical set of scenarios, and numeric results are obtained for them. The results indicate that to cover the serving area with the mentioned number of VNOs, subscribers, and SLAs, at least 11 OFDMA cells are required. The reference scenario assumes 16 cells, where the total network capacity is estimated as 1,800 Mbps. By reducing the number of cells to five, the total network capacity shrinks to almost half of its initial value (i.e. 48% of reference case). The changes in network capacity mostly influence the VNO with best-effort services, while the other types of VNOs suffer relatively less from the reduction of resources. Among guaranteed services, the violation starts from the service(s) with the lowest violation weight, which in our case study are background ones. The numeric results justify that the model is able to prioritise

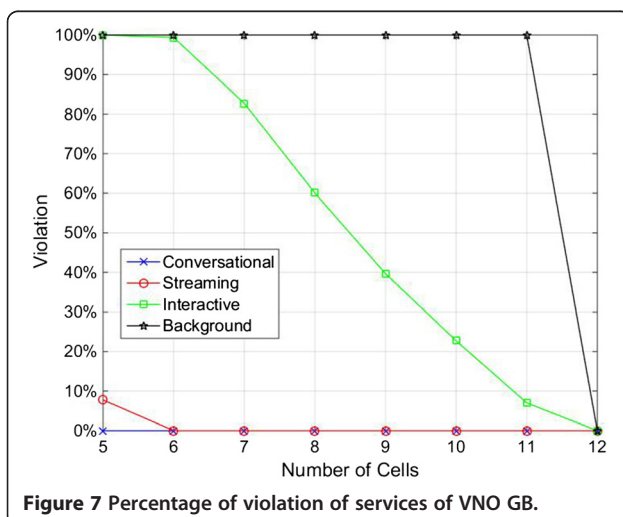


Figure 7 Percentage of violation of services of VNO GB.

the service according to their serving and violation weights. In the worst case studied in this paper, background and interactive services are totally shutdown while conversational ones experience no violation, as expected.

Competing interests

The authors declare that they have no competing interests.

Acknowledgements

The research leading to these results was partially funded by the European Union's Seventh Framework Programme Mobile Cloud Networking project (FP7-ICT-318109).

Received: 26 August 2014 Accepted: 10 February 2015

Published online: 12 March 2015

References

- Cisco Systems, *Global Mobile Data Traffic Forecast Update, 2012 - 2017* (Visual Network Index White Paper, Cisco Systems, Palo Alto, CA, USA, 2013). http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white_paper_c11-520862.pdf
- H Guan, T Kolding, P Merz, Discovery of cloud-RAN, in *Proc. of NSN cloud-RAN workshop, Beijing, China*, 2010
- D Sabella, P Rost, S Yingli, E Pateromichelakis, U Salim, P Guitton-Ouhamou, M Di Girolamo, G Giuliani, RAN as a service: challenges of designing a flexible RAN architecture in a cloud-based heterogeneous mobile network, in *Proc. of FuNeMS2013 - 22nd future networks and mobile summit, Lisbon, Portugal*, 2013
- R Yrjö, D Rushil, Cloud computing in mobile networks - case MVNO, in *Proc. of ICIN'11 - 15th international conference on intelligence in next generation networks, Berlin, Germany*, 2011
- A Klein, C Mannweiler, J Schneider, HD Schotten, Access schemes for mobile cloud computing, in *Proc. of MDM'10 - 11th international conference on mobile data management, Kansas City, MI, USA*, 2010
- C Mobile, *C-RAN - road towards green radio access network* (China, White Paper, China Mobile Research Institute, Shanghai, 2011) (http://labs.chinamobile.com/cran/wp-content/uploads/CRAN_white_paper_v2_5_EN.pdf)
- M Chiosi, D Clarke, P Willis, A Reid, J Feger, M Bugenhagen, W Khan, M Fargano, C Cui, H Deng, J Benitez, U Michel, H Damker, K Ogaki, T Matsuzaki, *Network function virtualisation: an introduction, benefits, enabler, challenges, and call for action* (France, White Paper, European Telecommunications Standards Institute, Sophia-Antipolis, 2012) (http://portal.etsi.org/NFV/NFV_White_Paper.pdf)
- Y Zaki, Z Liang, C Goerg, A Timm-Giel, LTE wireless virtualization and spectrum management, in *Proc. of WMNC'2010 - 3rd joint IFIP wireless and mobile networking conference, Budapest, Hungary*, 2010
- L Caeiro, FD Cardoso, LM Correia, Adaptive allocation of virtual radio resources over heterogeneous wireless networks, in *Proc. of EW'2012 - 18th European wireless conference, Poznan, Poland*, 2012
- Z Liang, L Ming, Y Zaki, A Timm-Giel, C Gorg, LTE virtualization: from theoretical gain to practical solution, in *Proc. of 23rd International Teletraffic Congress, San Francisco, CA, USA*, 2011
- X Costa-Perez, J Swetina, G Tao, R Mahindra, S Rangarajan, Radio access network virtualization for future mobile carrier networks. *IEEE Commun Mag* **51**(7), 27–35 (2013)
- T Metsch, P Gray (eds.), *"Infrastructure management foundations - specifications & design for mobile cloud framework"*, 2013. Deliverable D3.1, mobile cloud networking project, (<http://www.mobile-cloud-networking.eu>)
- S Khatibi, LM Correia, Modelling of virtual radio resource management for cellular heterogeneous access networks, in *Proc. of PIMRC'14 - IEEE 25th annual international symposium on personal, indoor, and mobile radio communications, Washington, DC, USA*, 2014
- A Khan, A Zugenmaier, D Jurca, W Kellerer, Network virtualization: a hypervisor for the Internet? *IEEE Commun Mag* **50**(1), 136–143 (2012)
- B Haberland, F Derakhshan, H Grob-Lipski, R Klotsche, W Rehm, P Schefczik, M Soellner, Radio base stations in the cloud. *Bell Labs Tech J* **18**(1), 129–152 (2013)
- J Pérez-Romero, X Gelabert, O Sallent, Radio resource management for heterogeneous wireless access, in *Heterogeneous wireless access networks*, ed. by E Hossain (USA, Springer, New York, NY, 2009)
- HS Dhillon, RK Ganti, F Baccelli, JG Andrews, Modelling and analysis of K-tier downlink heterogeneous cellular networks. *IEEE J Sele Areas in Comm* **30**(3), 550–560 (2012)
- Jacinto, N. M. d. S., "Performance gains evaluation from UMTS/HSPA+ to LTE at the radio network level", Master of Science, Department of Electrical and Computer Engineering, Instituto Superior Técnico, University of Lisboa, Lisbon, Portugal, 2009 (http://grow.inov.pt/wp-content/uploads/2014/01/2009_NunoJacinto.pdf).
- A Papoulis, SU Pillai, *Probability, random variables, and stochastic processes* (McGraw-Hill, New York, NY, USA, 2002)
- M Buehrer, RM Buehrer, *Code division multiple access (CDMA)* (Morgan & Claypool Publishers, San Rafael, CA, USA, 2006)
- MATLAB and Statistics Toolbox Release 2012b, The MathWorks, Inc., Natick, Massachusetts, United States, <http://www.mathworks.com>, Accessed Feb, 2014.
- Y Zhang, *Solving large-scale linear programs by interior-point methods under the MATLAB environment, Technical Report TR96-01* (University of Maryland, Baltimore, MD, USA, 1995)
- S Mehrotra, On the implementation of a primal-dual interior point method. *SIAM J Optim* **2**, 575–601 (1992)
- C Systems, *The Zettabyte era - trends and analysis* (Visual Network Index White Paper, Cisco Systems, Palo Alto, CA, USA, 2013) (http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/VNI_Hyperconnectivity_WP.pdf)
- C Systems, CVN Index, *Global mobile data traffic forecast update, 2011-2016, from visual network index (VNI) white paper* (Cisco Systems, California, USA, 2012)
- Scalable and adaptive Internet solutions (SAIL) project, (<http://www.sail-project.eu/>), 2015.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► springeropen.com