



Contents lists available at ScienceDirect

Computer Communications

journal homepage: www.elsevier.com/locate/comcom

QoE-driven dynamic management proposals for 3G VoIP services

Jose-Oscar Fajardo^{a,*}, Fidel Liberal^a, Is-Haka Mkwawa^b, Lingfen Sun^b, Harilaos Koumaras^c^a NQaS Research Group, University of the Basque Country, ETSI Bilbao, Alameda Urquijo s/n, 48013 Bilbao, Spain^b Centre for Signal Processing and Multimedia Communication, School of Computing, Communications and Electronics, University of Plymouth, Plymouth PL4 8AA, UK^c NCSR DEMOKRITOS, Institute of Informatics and Telecommunications, Agia Paraskevi, 60228 Athens, Greece

ARTICLE INFO

Article history:
Available online xxx

Keywords:
QoE
E-model
VoIP UMTS
Cross-layer adaptation

ABSTRACT

Most of the currently available adaptation solutions of VoIP over UMTS are based on the modification of service parameters as the only available reaction against any detected service degradation. On the contrary, in this paper we propose a combined approach where service-level adaptation is considered first and, provided that no suitable parameter combination is capable of providing enough QoE, a change of network state will be suggested. In order to do so we analyze the performance of the end-to-end (e2e) performance metrics in this convergent scenario, the root causes of possible degradations and, finally, the combined effects of the different network segments and their impact on the user perceived QoE. We show the map of best performing VoIP configurations for every state of the network segments. Furthermore, considering each of these configurations, we analyze the acceptability of the service or the convenience of trying to modify the network state. Finally, a lightweight implementation based on simple network state estimation and decision heuristics is proposed and validated in terms of accuracy and responsiveness.

© 2010 Elsevier B.V. All rights reserved.

1. Introduction

One of the most relevant examples of converged networks is the Universal Mobile Telecommunications System (UMTS). From its design, the system was conceived as a multipurpose network capable of simultaneously supporting multimedia services over the Circuit-Switched connections, and data connections over a Packet-Switched approach. However, the evolution of the capabilities offered by the current mobile Internet accesses, jointly to the generalization of the multimedia services in the Internet, lead the trend towards a converged network where both data and multimedia services are supported over the mobile data plane.

Within this new world of business opportunities some specific services can be taken as reference, namely the mobile Voice over IP (VoIP) and the different modalities of mobile video provisions. One of the drawbacks in the provision of this kind of multimedia mobile services is the difficulty in assuring some level of perceived Quality of Service (QoS) which shows an impact on the global Quality of Experience (QoE) levels. To the typical network performance variability found in the Internet, we have to add the high variability associated to the performance of the wireless access networks.

The key challenge to overtake this network variability is the concept of adaptation, either at service or network level, or even better in a combined cross-layer basis.

Nowadays, most of the network operators do not rely on the deployment of an accurate network management system in order to offer QoE-aware mobile multimedia services over data connections. It seems that the complexity associated to this kind of solutions does not convey the migration in terms of business revenue.

On the contrary, most of the currently available adaptation solutions are based on the modification of service parameters as a reaction of the monitored service performance. However, the decision making engine under this adaptations is usually quite trivial, and not QoE-aware. Typical adaptation actions include the decrease of the required bitrate when network packet losses are detected, or the modification of the play out buffer upon variations in the network jitter.

This kind of default adaptation actions are not always the most optimal from the QoE perspective. And this fact is even more relevant in a mobile Internet service, where the performance degradations may be caused by a series of heterogeneous reasons that cannot be mitigated by the same reactions. Thus, we address the need for a suitable analysis of the possible causes of degradations, before any adaptation action can be proposed.

This paper goes in depth into the analysis of the performance of the end-to-end (e2e) metrics, and their impact on the user perceived quality of service for mobile VoIP services as a dimension of the QoE. From the study of the root causes of degradation, we

* Corresponding author. Tel.: +34 94 6017361; fax: +34 94 6014259.

E-mail addresses: joseoscar.fajardo@ehu.es (J.-O. Fajardo), fidel.liberal@ehu.es (F. Liberal), is-haka.mkwawa@plymouth.ac.uk (I.-H. Mkwawa), L.sun@plymouth.ac.uk (L. Sun), koumaras@iit.demokritos.gr (H. Koumaras).

will analyze the combined effects of the different network segments and the expected QoE for different network states. These results will be the basis to determine the most suitable configuration of the service chain for each combined network state, and will allow us to propose the optimal adaptation (or set of cross-layer adaptations) based on the available feedback information. Specifically, the contribution of this paper is threefold.

- First, a service-level adaptation map is proposed for the dynamic management of VoIP over 3G services. Based on the knowledge of the performance associated to each network segment, we show the best performing VoIP configuration as a distributed solution, taking into account the adaptation capabilities at both service endpoints. Thus, each proposed VoIP configuration results on the optimal performance for the experienced network performance in a coordinated way, instead of making individual- and likely suboptimal-decisions.
- Secondly, based on the QoE level of each selected configuration, we analyze the acceptability of the service since even the best service-level option may result on poor quality. In case of unacceptable service grades, the decision map provides the information concerning the quality level expected in the adjacent network states. Therefore, this information can be used to detect the minimum network-level adaptations required. Since the new network state may involve a new optimal VoIP configuration, the outcome of the decision making results on a QoE-driven cross-layer adaptation procedure.
- Finally, a lightweight implementation of the system is proposed. It may be difficult to obtain real-time feedback about the performance of the access and the core networks. Hence, we propose how these network states can be inferred from the service endpoints. Likely being mobile devices with limited capacity, we propose a simple estimation and decision procedure, and validate this approach in terms of accuracy and responsiveness.

The remainder of the paper is structured as follows: In Section 2 the motivation of this work is analyzed, visiting the fundamentals and trends of the related dynamic service adaptation approaches. Section 3 focuses on the analysis of how to estimate the QoE level from the combined service and network status, and proposes a mechanism for the selection of the best scoring configuration. Section 4 addresses the specific topic of QoE-driven management of VoIP over 3G services. After the specification of the service characteristics, we show the most optimal configurations at service level first, and from a cross-layer perspective later. Once the decision making plane is analyzed, Section 5 deals with the implementation issues of the proposed system. Specifically, this section studies how the network information can be inferred from the e2e metrics, enabling the implementation of the logic at the service endpoints. Finally, Section 6 focuses on the validation of the whole decision system and Section 7 gathers the main conclusions to the paper.

2. QoE-aware adaptation mechanisms

When addressing QoE management in complex environments such as the one considered in this work, two main trends are considered: network-aware services and service-aware networks. Both approaches highlight the need for deploying combined mechanisms that allow optimizing the content generation and service provisioning procedures in a personalized way in order to assure an enriched QoE for end users. However, the service entity where the adaptation is carried out usually differs in each case. While the former concerns the fine configuration and adaptation of service parameters in function of the expected or experienced

network performance, the latter focuses on the best usage of the network resources taking into account the content characteristics.

The concept of VoIP service-awareness in the management of converged networks has been the objective of many studies for long now. For example, in [1] it was already overviewed the main issues in the heterogeneous management of resources between the UMTS access network and the interconnecting core network. Concerning the network-aware service, most of the work has been devoted to the voice codec switching under degraded network conditions [2]. In this sense, the Adaptive Multi-Rate (AMR) codec, standardized by the 3GPP in 1998, offers two main features that makes it the predominant candidate for UMTS Packet-Switched services: accurate speech quality levels at low data rates, and the capability to switching the codec mode almost in real-time, in a per-frame basis (each 20 ms) [3]. The narrowband AMR codec can work at eight codec modes, with different target encoding bitrates from 12.2 kbps to 4.75 kbps. Each lower codec mode provides a lower listening quality level in ideal service conditions. Yet, if network impairments occur, the codec may change its configuration to a lower codec mode in order to prevent deeper degradations.

The analysis of the main sources of degradation provides the required knowledge for proposing the most suitable adaptation actions. Although VoIP QoE modeling will be revisited in Section 3 it is commonly known that the most important network factors are the end-to-end delay of the communication and the experienced packet loss. In [4], authors provide a comprehensive analysis of the different contributions to the delay and loss ratio in a UMTS-to-PSTN scenario. The end-to-end delay can be estimated as a composition of service-layer delays, and delays introduced by both the UMTS Access Network (AN) and the Core Network (CN) – see the following equation:

$$d_{e2e} = d_{VoIP} + d_{AN} + d_{CN} \quad (1)$$

The delays introduced by the VoIP service (d_{VoIP}) are caused by the codification and decodification delays (i.e., d_{cod} and d_{decod}), the delay introduced by the packetization scheme (d_{pack}) and the delay introduced by the dejittering buffer (d_{jit}). The delays introduced by the network can be divided into two main factors. The delay introduced in the UTRAN (d_{AN}) includes the transmission delay and the additional delay due to Radio Link Control (RLC) functions. This RLC delay is made up of the delay due to the slotted access to the physical medium and the delays introduced by the RLC Protocol Data Unit (PDU) recovery mechanism. The latter highly depends on the UTRAN loss ratio, which determines the number of RLC PDUs to be recovered, and on the burstiness of the losses, which may influence the time needed for recovering a erroneous SDU. Finally, the delays introduced by the Core Network (d_{CN}) are mainly due to the transmission effects and the possible additional delays in the queues due to congested nodes.

Likewise, the total packet loss ratio (ρ_{e2e}) can be determined by the individual contributions of each packet loss ratio due to the different sources of losses.

$$\rho_{e2e} = f(\rho_{AN}, \rho_{CN}, \rho_{jit}) \quad (2)$$

ρ_{AN} is the contribution to the packet loss ratio caused by packets lost in the UTRAN. Considering the UTRAN working on RLC AM, these packet losses are caused by the impossibility to locally recover the lost RLC PDUs, e.g., due to the expiration of the SDU Discard Timer. ρ_{CN} is the packet loss ratio at the core routers, mainly due to possible congestion situations. Finally, ρ_{jit} is the packet loss ratio experienced in the dejittering buffer, which is caused due to voice frames arriving later than its play out time. Thus, the value of ρ_{jit} is highly impacted by the statistical characteristics of the end-to-end delay and the dejittering process.

According to the definition of these e2e metrics, the VoIP over UMTS scenario considered in this study may include different possible dynamic adaptation actions at service level:

- *Change in the AMR mode.* Upgrade to an upper AMR mode means a better audio quality at the cost of higher bitrate requirements. Downgrade to a lower AMR mode would result on a worse audio quality on source, but it lowers the bitrate requirements and, thus, the media flow may avoid degradations through the network.
- *Change in the packetization scheme.* Including more than one voice frame per RTP packet decreases the required bitrate as well. However, it introduces additional delays to the communication and higher packet sizes make the packet loss probability increase.
- *Change in the dejittering buffer size at the destination end user.* An increase in the dejittering buffer size would likely lead to lower voice frame losses within the buffer itself. However, it introduces an extra delay into the play out function. Thus, it is a trade-off to determine the most suitable value for each case, in function of the mean value and the variability of the end-to-end delay.

The simplest default adaptation action is limited to decreasing the voice quality until the effect of network impairments is mitigated. Once the network recovers its previous performance level, the codec may be upgraded again to improve the encoding audio quality. These transitions are driven by a hysteresis cycle, as shown in [5] for an IMS platform, and the result offers a significant quality enhancement over the case of no adaptations.

Many enhancements to this trivial procedure have addressed the selection of the best performing AMR mode for different UMTS radio conditions. Yet, as explained in [6], the selection of the AMR mode in VoIP services must cope with the best configuration for all the possible combinations of access network and core network states. This paper proposes a network-initiated selection of the AMR mode, taking into account for the decision the feedback information about the performance of the different network segments. As a drawback, the adaptation actions considered in this paper are limited to the change of the AMR mode, while other VoIP configuration parameters are set up to constant values. For instance, the packetization is kept at one single voice frame per IP packet, which may not be the optimal configuration when the total delay for the VoIP session takes tolerable values.

Likewise, Ref. [7] focuses on the analysis of the optimum dejittering buffer size for different radio quality states, finding an optimal buffer size of 60 ms for low values of the Block Error Ratio (BLER). Several limitations are found in these results. On one hand, the studied buffer size value is increased from 60 ms to 180 ms, without specifying intermediate values. Moreover, the set of adaptation actions is also limited to the modification of the buffer size, and the delay introduced by the core network is not considered as well.

The contribution of all the network segments to the e2e delay and loss ratio is well-analyzed in [4]. The proposed procedure of analysis is useful to study the combined contributions to the e2e performance metrics, and to their impact on the QoE based on the E-model. As an example, the procedure is applied to detect the most suitable combinations of packet size and dejittering buffer size for different RLC configurations. However, in this proposal the radio link quality is not taken as feedback information for the AMR mode or codec switching process.

Beyond all these particular approaches, the most spread implementation of the e2e adaptation capabilities for VoIP nowadays is likely the Skype application. Although being a private implementation, several works have tried to understand its operation from

experimental conditions. Thus, authors of [8] infer the QoE-oriented reactions that Skype performs upon network degradations. Three types of adaptation actions are detected: (i) when network losses appear, the sent packet sizes increase at the same packet rate, so it seems that some kind of redundancy is used; (ii) if the degradation persists and reaches about 20% of packet loss, the voice codec is switched to a lower bitrate; (iii) if the delay between the two endpoints becomes unacceptable, the data flow is relayed over a third-party node trying to avoid a possible congested link in the Internet.

Therefore, Skype seems to implement some kind of edge-based intelligence, trying to discern the source of network degradation and performing the most suitable adaptation in reaction, including service-level and network-level modifications of the VoIP provision chain. In the case of service-level adaptations, two main drawbacks can be stated from [8]. On one hand, the application does not react to bursty losses, but only to random losses. On the other one, the response time is about 1 min, which is a high value compared to the per-frame adaptation capability of the AMR codec. With regard to the network-level adaptation, although the response time is reduced to 15 s, it seems to start only for very high delay values of about 4 s. Finally, no information is available about the possible joint configuration of the dejittering buffer based on the experienced performance of the combined effects of the total delay and the ratio of packet losses.

Recently, several works have been focused on the topic of cross-layer optimization of the provision chain, aimed at maximizing the QoE of multimedia services over mobile accesses. For example, in [9,10] authors propose a greedy resource allocation algorithm to solve the problem of optimizing the performance of a HSDPA cell. These studies study how to optimize the provision of services when the resources of the access network are not enough to cope with the accumulated of users requirements. Thus, the radio resource allocation process is driven by QoE-related utility functions, instead of the traditional throughput-maximization approach. For the evaluation of the different solutions, the variable argument includes both service-level parameters and network-level parameters. Concerning the VoIP service, the evaluation of the service for each individual user is based on which of the considered codecs offers a better performance for the possible network states, taking into consideration the rate and the packet error probability that can be achieved at each network state evaluated. Authors propose a greedy algorithm for solving the optimization problem, and show how the QoE-driven optimization of the cell offers enhanced results in multi-user contexts.

However, these kinds of solutions are difficult to implement in a case study such the proposed in this paper. As cited before, the converged approach entails that the end-to-end network performance metrics are determined by the contributions of the different network segments. Considering random patterns for the delay in both AN and CN segments, the random variable associated to the total delay should be modeled as the convolution of the random delay variables of each segment. Thus, the modeling of the end-to-end network performance becomes quite complex. Moreover, the resulting network model could be too complex to be included in an optimization algorithm, even more if we include other service-level parameters as the packetization scheme or the buffer size into the equation.

As a result, the optimization approach considered in this study addresses the problem from a different perspective. First, the combined impact of all the considered service configurations and network states are evaluated by means of intensive simulations. Later, the *a priori* knowledge will be used for the cross-layer adaptation decision making process.

This paper tries to overcome this set of partial solutions, providing a comprehensive study of the most optimal 3G UMTS VoIP

service adaptation procedures in different combined AN/CN states. Fig. 1 illustrates the basics of the adaptation procedure that is proposed. Each horizontal plane represents a combination of AN and CN performance states. For each coordinate within the plane, the vertical axis represents all the possible combinations at service level. Each combination in the 3D space will take an associated estimated QoE value, which will lead the decision making process for the adaptation system.

Therefore, this proposal tries to cope with the following objectives:

- The adaptation procedure will implement an edge-based intelligence, being the endpoints capable of monitoring the network states and proposing the most suitable adaptation actions.
- The detection of the source of the network degradations (or the combination of effects) will cope with the requirement of real-time reactions.
- When a new combined AN/CN state is detected, the automated system should infer the optimal VoIP configuration for the new network state. This means that the system will be placed at the best scoring position in the vertical axis. For the new VoIP configuration selected, the system will be able to estimate if the expected QoE fulfils the required quality.
- In that case, if the VoIP configuration differs from the current one, the system will launch the procedure to take the required adaptation action, which could affect any or both of the endpoints: the sender endpoint, if the AMR mode or the packetization scheme shall be modified, and the receiver endpoint, if the de-jittering buffer size shall be changed.
- Otherwise, the only possible solution is to move to another AN/CN state where the service provision is able to comply with the quality requirements. This action is represented as network adaptation in Fig. 1, and it will likely entail a possible change in the network or in the Quality of Service (QoS) classification of the involved traffic. If the new position in the horizontal plane requires the modification of the VoIP session to the corresponding best configuration, this will represent a cross-layer adaptation.

The former is basically a network-driven service adaptation approach, and can be implemented from the endpoints of the communication. The latter becomes a more complex cross-layer service and network adaptation, and its real-world implementation involves the capability to request modifications in the utilization of network resources. Thus, this approach could be interesting from a network operator perspective. When considering multi-user scenarios, a network adaptation may impact the quality perceived by the rest of the users with shared network resources, so the

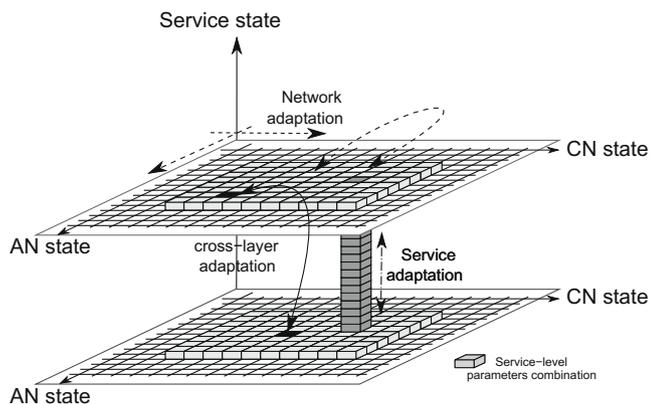


Fig. 1. Network/Service + cross-layer adaptation approach.

decision making could include some kind of cross-layer optimization specific to the considered network segment. However, this feature is not considered in the scope of this paper.

3. Analysis of QoE of VoIP services

The predominant method for assessing the PQoS of VoIP services is the E-model, defined by the ITU-T in [11]. This model allows computing the value of a rating factor R , which provides an evaluation of the communication impairment as shown in the following equation:

$$R = R_0 - I_s - I_d - I_{e-eff} + A \quad (3)$$

The R score computed with default values as proposed in [11] is 93.2, which corresponds to an estimated Mean Opinion Score (MOS) score of 4.41. The additional impairment due to the network performance is caused by two main factors. I_d is defined as the impairment due to the total mouth-to-ear delay, while I_{e-eff} is the impairment factor due to the combined effect of the codec and the packet losses. The formulation for computing the I_{e-eff} factor introduced in the 03/2005 version of the E-model defines a complex relationship between the packet loss ratio (P_{pl}), the characteristics of the losses concerning its burstiness ($BurstR$), the codification impairment (I_e) and the codec robustness to packet losses (B_{pl}).

$$I_{e-eff} = I_e + (95 - I_e) \cdot \frac{P_{pl}}{\frac{P_{pl}}{BurstR} + B_{pl}} \quad (4)$$

One of the key advantages of these expressions is that they provide a way to estimate the expected quality of VoIP services under specific service and network conditions. For the aims of this study, we focus on the analysis of the quality levels achievable for AMR-based VoIP services. A comprehensive study of the application of the E-model to AMR-based VoIP services over 3G UMTS connections is presented in [12]. Unlike I_d , which is codec-independent, the value of I_{e-eff} is determined by both the performance metrics (P_{pl} and $BurstR$) and the codec-dependent coefficients (I_e and B_{pl}). Therefore, in order to carry out a comparative analysis of the performance for different VoIP configurations, the values of P_{pl} and $BurstR$ are required for the considered codecs.

The AMR-12.20 mode behaves as the GSM-EFR codec, which is included in the specifications provided in Appendix I/G.113 [13]. Therefore, the considered values are $I_e = 5$ and $B_{pl} = 10$. Likewise, the AMR-7.40 codec mode is compatible with the IS-641 codec, which takes an I_e value of 10 with no B_{pl} value specified. Following the E-model, considering a default value and delay of 150 ms and under loss-free conditions, the maximum achievable MOS scores are 4.288 for AMR-12.2 and 4.133 for AMR-7.4.

However, the values of these parameters are not available for the whole set of the AMR modes. From the scientific bibliography, one of the most comprehensive studies regarding AMR subjective quality ratings is provided in [14], where a method for computing the I_{e-eff} factor in a logarithmic form (5) is proposed.

$$I_{e-eff} = a \cdot \ln(1 + b \cdot \rho) + c \quad (5)$$

In the expression (5) ρ is the voice packets loss rate while a , b and c are constants dependent on AMR mode as shown in Table 1.

Compared to the original I_{e-eff} expression, the c coefficient would be the corresponding to the I_e factor, while the impairment and codec resiliency to packet losses become determined by the logarithmic factor. Thus, this expression will be used in the estimation of the MOS values for the different network conditions.

Table 1Fitting parameters for I_e vs. packet loss for different codecs (based on PESQ-LQO).

Par.	AMR-122 (12.2 kb/s)	AMR-102 (10.2 kb/s)	AMR-795 (7.95 kb/s)	AMR-74 (7.4 kb/s)	AMR-67 (6.7 kb/s)	AMR-59 (5.9 kb/s)	AMR-515 (5.15 kb/s)	AMR-475 (4.75 kb/s)
a	22.98	21.14	22.80	22.63	22.86	23.41	25.83	26.46
b	0.305	0.362	0.220	0.211	0.180	0.148	0.100	0.088
c	10.07	13.23	19.50	20.76	23.79	27.36	30.45	32.42
R^2	0.9997	0.9999	0.9998	0.9999	0.9999	0.9999	0.9999	0.9998

3.1. QoE estimations for the case study of converged networks

The previous expressions allow us to develop comparative analyses of the expected QoE levels based on the computation of the MOS scores. Considered a specific network performance status, the loss and delay patterns can be obtained. Jointly with the buffer configuration, these metrics determine the total mouth-to-ear delay and the total loss ratio of voice frames. The expression (5) allows us to follow the impairment produced by the combined effect of the service configuration and the network performance. And finally, the evolution of the MOS scores will determine the user experienced quality. As well, the inclusion of the *BurstR* parameter in (4) allows estimating the expected MOS scores for different packet loss conditions. The *BurstR* parameter reflects the short-term dependency of packet losses, as detailed in [15]. This is especially useful when considering radio-based mobile communications, due to the bursty nature of losses.

However, we face two main problems for the direct application of these expressions for the considered converged network scenario. On one hand, the *BurstR* parameter is related to the IPI-layer packet loss pattern. As we describe in Section 4.1.2, our case study considers the performance of the AN to be based on the RLC-layer error model. Although this error model is based on a two-state Markov model, it cannot be directly mapped to the IP-layer packet loss model when RLC retransmissions are enabled since this relationship depends on several parameters, such as the error rate, the available bitrate for RLC retransmissions, the retransmissions timer or even the packets size.

Additionally, the packet loss pattern in our case study is a result of the contributions of the different network segments. Even under no network losses, the final impact on the expected quality is a function of the different contributions to the total delay, which determines the losses at the de jittering buffer. If we consider random delay contributions, the total network delay can be obtained from the convolution of the individual contributions. As a result, we need to run long voice simulations in order to capture all the possible effects of the combined contributions of the random variables for a single combined AN/CN state. So, a long-term evaluation of the quality is required.

Fig. 2 shows a comparison between the evolution of the MOS scores for two different VoIP sessions, with different service configurations, which are being provided under the same network conditions.

From the inspection of these traces, we can state the difficulty of selecting the optimal result. The high variability of the MOS scores over time found in this case illustrates that the mean MOS values for the whole trace is worthless for representing the user experienced QoS by itself.

In order to evaluate the global quality of a long sequence, where the network conditions vary over the time, an integral approach shall be used. In [15], different approaches intended to solve this problem are analyzed and compared. If we consider periods of dependant losses under macroscopic evolutions of the packet loss, four approaches are considered in order to obtain the integral quality from the time-averaged quality assessments for each period. A weighting process is adopted for this study,

since it shows better results for the integral quality estimation under severe degradations.

However, the analyzed weighting method is based on two aspects: the strength of the degradation over a time interval and its relative position within the complete stimulus. The latter is also considered in this study, assuming that users weight strong degradations more critically. However, the former is excluded from this study. It is well known that the recency effect may likely modify the integral quality assessment of an individual. Yet, since the aim of this study is to detect the expected QoE level associated to each considered combined AN/CN state, we need to eliminate the effect of the random appearance of severe network impairments due to combined contributions.

For the aims of this study, we introduce a new objective QoE indicator as the MOS_{factor} , which is aimed at including the previously discussed features of the long-term evaluation of quality.

It is generally accepted that it is preferable to keep a lower – but still acceptable – mean MOS value if these results on a lower number of annoying degradations over the whole trace. For example, in [15] it is shown that users' reactions to quality improvements are slower than their reactions to quality degradations. Thus, the percentage of MOS scores at several specific points arguably provides a better estimation of the global QoE than the average value itself. In this sense, the Experimental Cumulative Distribution Function (ECDF) is considered to be useful to depict the general behaviour of a sample variable over time. For instance, Fig. 3 shows the ECDF of the MOS scores found for a series of different VoIP configurations, all them under the same network conditions.

This kind of results can be useful for analyzing if the target QoE requirements are fulfilled. For instance, a good QoE level could require that the computed MOS scores are over 3.5 for the 99% of the time. Thus, depending on the policy implemented for each user, this analysis could provide the feedback about the grade of QoE achieved.

Based on a typical interpretation of the MOS scale, we define the following significant intervals in comparison to the Public-Switched Telephone Network (PSTN):

- $I_1 = MOS > 3.5$: most users perceive a voice quality similar to the PSTN quality.
- $I_2 = 3.1 < MOS < 3.5$: under PSTN service, but acceptable quality.
- $I_3 = 2.5 < MOS < 3.1$: quite annoying quality for most users.
- $I_4 = MOS < 2.5$: quality level becomes unacceptable.

According to the definition of these intervals, we establish the following QoE requirements for this study: the number of MOS samples in the I_4 interval should be kept under the 1% of the time, while the requirement for the I_3 interval is set to the 10% of the time. These two conditions are established as mandatory, so their significance is very high. If the results for these two lower intervals are similar, the values within I_2 and I_1 become the selection criteria for the best performing configuration.

Thus, the MOS_{factor} is finally defined as:

$$MOS_{factor} = 100\% \{ECDF_{I_4}\} + 10\% \{ECDF_{I_3}\} + 1\% \{ECDF_{I_2}\} + 0.1\% \{ECDF_{I_1}\} \quad (6)$$

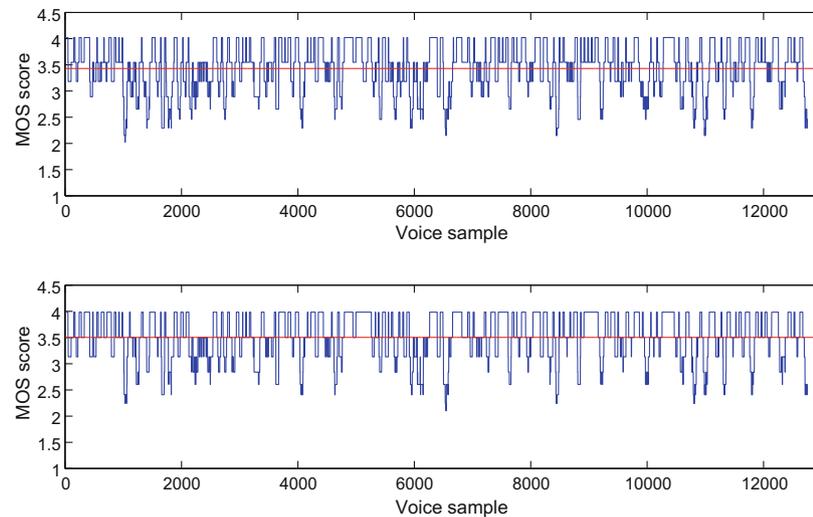


Fig. 2. Comparison of the MOS evolution of two VoIP sessions.

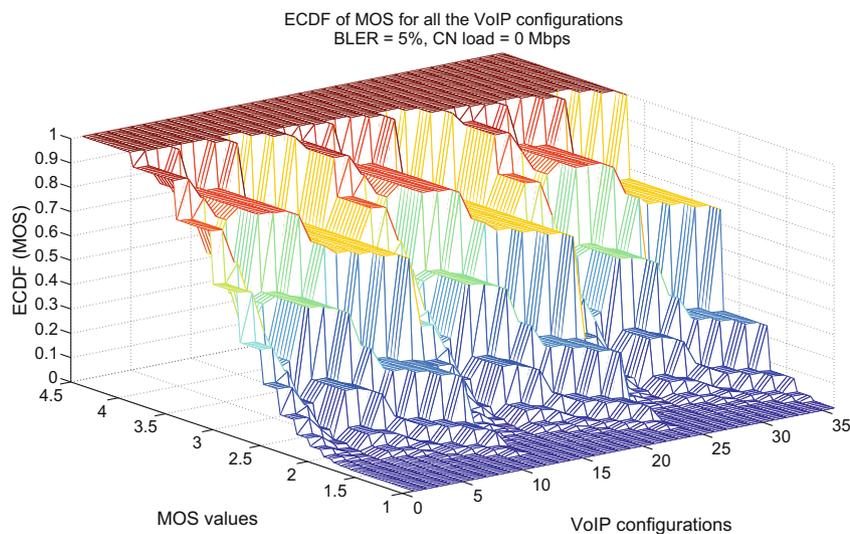


Fig. 3. ECDF of MOS scores for different VoIP configurations.

As a result, the best scoring configuration is that with the lower MOS factor. Fig. 4 illustrates the result of the proposed weighting process. As can be observed, the MOS scores for the two cases represented seem very similar, and they would produce barely any difference in the typical statistics. Yet, the second case provides a lower number of severe degradations (below the 2.5 threshold) and so is reflected after the weighting process. Therefore, the proposed approach proves to be efficient in the detection of this kind of behaviours.

4. VoIP services over 3G UMTS networks

4.1. Analysis of the service provision

The considered VoIP service provision is based on a typical best-effort Internet connection, without specific QoS procedures that would offer network performance guarantees to the multimedia flow.

4.1.1. VoIP characteristics

In this paper, we focus on the AMR-based VoIP over 3G services [16]. This codec offers a good quality at low bitrate demands and a

real-time adaptation capability, and its performance AMR in 3G channels has been proven in [3]. The AMR codec may work on eight modes at different bitrates. The voice frame is always 20 ms length, with frame sizes from 95 to 244 bits to fit the target bitrate.

The total header overhead is 40 bytes per voice frame (12 bytes RTP, 8 bytes UDP and 20 bytes IP). Header compression is not considered in this study, since it has been proved that it may induce to more severe packet losses under the same radio link conditions [17]. The main alternative to reduce the protocol overhead is by modifying the packetization, including a higher number of voice frames per IP packet. Yet, due to latency considerations in this kind of interactive services, the number of frames included within each RTP packet cannot be increased too much. Thus, the optimal packetization will depend on the rest of contributions to the total delay.

Finally, the impact of the network performance on the perceived service level depends on the de-jittering buffer size. This buffer is introduced in the destination endpoint between the receiving and playing functions in order to store a determined number of voice frames. Thus, voice frames are delayed before being played out, but the application results more resilient to delay variations. As a result, the effect of the network jitter is reflected in the performance of the de-jittering buffer. If a delayed VoIP packet arrives at

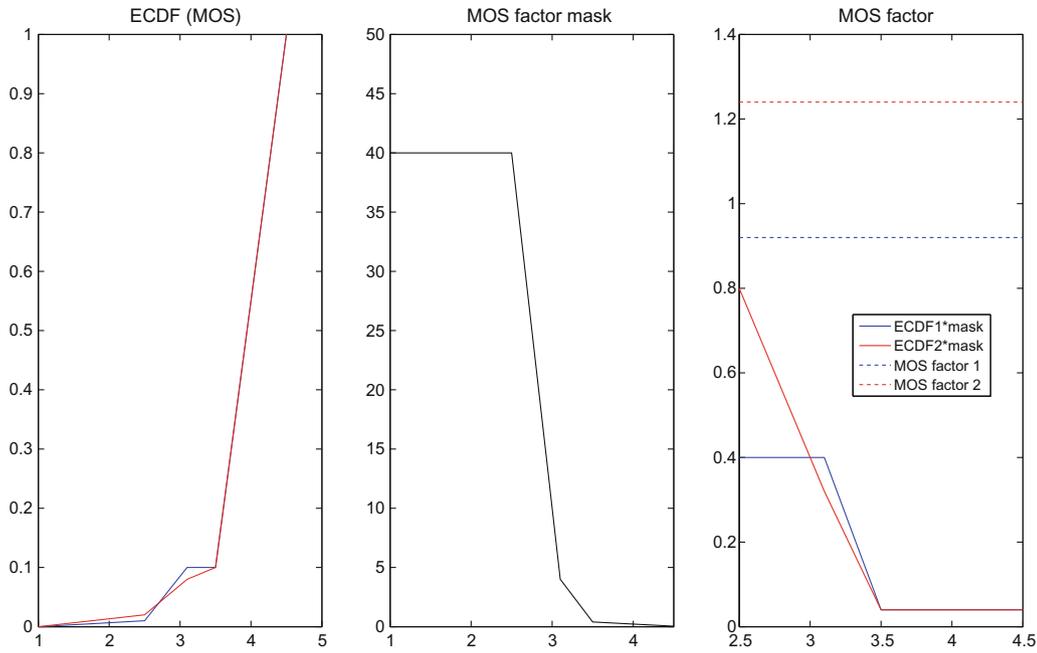


Fig. 4. MOS_{factor} weighting process and result.

the destination while the subsequent voice frame is still in the dejittering buffer, it could be provided to the playing function on-time. Otherwise, the packet will be computed as lost in the dejittering buffer.

4.1.2. UTRAN model

With regards to the UTRAN configuration, the VoIP service is supported over a Background Packet Data Protocol (PDP) Context with a typical mobile wide area configuration as defined in 3GPP TR 25.993 [18] for the “Interactive or Background/UL:64 DL:384 kbps/PS RAB”. The transmission channel supports maximum bitrates of 384 kbps in the downlink and 64 kbps in the uplink over a Dedicated Channel (DCH). If no header compression is considered, the maximum required bitrate for AMR over IP services reaches 28.8 kbps, so the available bandwidth seems high enough for AMR services. Table 2 shows the most relevant parameters considered for the UTRAN.

Due to this configuration, when an RLC PDU is lost, the RNC tries to retransmit this PDU. When all the PDUs of a RLC SDU are correctly received, the UE sends it to the upper layer regardless the status of the previous RLC SDUs. If a retransmitted RLC PDU is once again lost, the RNC tries the retransmission until the SDU Discard Timer expires.

Table 2
UTRAN considered parameters.

UTRAN feature	Value
<i>RLC layer</i>	
Max. bitrate at RLC level	384 kbps
RLC PDU size	320 bits
RLC mode	Acknowledged Mode (AM)
Delivery order of SDUs	Not in-order delivery
Allowed Transport Format Set (TFS)	Six possible TFs: 0–1–2–4–8–12 TB/TBS
SDU discard mode	Timer based, discard timer = 500 ms
SDU concatenation	Enabled
<i>PHY layer</i>	
Transport Block (TB) size	336 bits
Transmission Time Interval (TTI)	10 ms
Transmission Channel (TrCH) type	Dedicated Channel (DCH)

Although the considered RLC AM service supports the recovery of radio errors in the UTRAN, the quality of the audio reception may be impacted in several ways. The local recoveries introduce additional delays, which may lead to frame losses in the application buffer. As well, the local recoveries are limited by a counter, so in severe radio degradations some frames may be actually lost in the UTRAN. Additionally, these recoveries increase the required bitrate in the radio channel, which in high usage ratios as in the uplink may introduce additional delays to the transmission of new voice frames. As a result of all these considerations, we can state the great relevance of the impact of the specific UMTS error conditions.

Due to the bursty nature of the link error pattern in the wireless medium, the implemented UMTS link layer model is based on a two-state Markov model. This model governs the performance of the UTRAN at RLC level, thus determining the correct reception of Transport Blocks (TB) at the user side. The relevant parameters for evaluating the impact of the link losses into the service performance are the Block Error Rate (BLER), which determines the ratio of Transport Blocks received with errors, and the Mean Burst Length (MBL), which introduces the grade of burstiness in the error pattern. The relation of these two parameters with the transition probabilities are given by the following equations:

$$P_{ee} = 1 - \frac{1}{MBL} \quad (7)$$

$$P_{ce} = \frac{(1 - P_{ee}) \cdot BLER}{1 - BLER} \quad (8)$$

The model parameters are configured based on the results presented in [19] from the monitoring of live 3G UMTS connections. Specifically, two characteristics have been considered:

- For mobile users, the radio errors can be aggregated at Transmission Time Interval (TTI)-level. This means that when errors occur in a TTI, most likely all the TBs within that TTI will be erroneous.
- The bursts of errors produced in the UTRAN follow a MBL of approximately 1.75. Thus, the probability that more than two or more TTIs are consecutively in error is about the 40%.

4.1.3. CN model

In what concerns to the core network, the transmission of data is modeled as a basic transport service. Provided that the transmission links in the core network do not experience physical problems, the overall network performance will be determined by the performance of the routing elements.

In the typical best-effort service, each node examines the next hop of the IP packet and then places the packet into the corresponding output interface. The variability of the network metrics is caused by the behaviour of these output queues. The latency in the transmission of a packet depends on its transmission time, as well as the transmission time of all the IP packets that are previously stored in the transmission queue, being the transmission time a function of the packet size and the transmission capacity. This effect determines the contribution to the total delay.

Additionally, the transmission queues are of finite size. Thus, if a burst of packets arrives at the node the queue could be congested-physically or logically by configuration-causing packet drops.

Therefore, the general performance will be determined by the combination of the node parameters – output capacity and transmission queue size – and the characteristics of the arriving traffic – the packet arrival rate and the packet size distribution. In this study, the combination of the output capacity and the packet arrival rate will be established by a unique configuration parameter, namely the traffic intensity. Specifically, we consider the case of a traffic trunk that is exclusively devoted to the transportation of all the data packets belonging to the UMTS traffic. Thus, we consider a traffic trunk of 2 Mbps which will work at different operating points depending on the evolution of the supported traffic intensity.

With regard to the packet size, the expected traffic is considered to be made up of the combination of packets belonging to different VoIP and mobile video flows.

- For VoIP flows, the packet sizes are kept at low values. For example, for the three higher AMR modes and for packetization schemes up to 3, the packet sizes fall in the range of 60 bytes to 131 bytes [12].
- For mobile video resolutions, the optimal value for the maximum packet size is quite linked to the optimal video slice size, which at the same time depends on the experienced BLER value, as explained in [20]. Different studies propose values of 110, 300 or 680 bytes as optimal values for different BLER conditions. For good radio reception conditions, values close to the MTU can be used, e.g., 1400 bytes of maximum slice size.

Experimental results show that, keeping a comparable number of simultaneous VoIP and video users with varying flow characteristics, the packet size distribution can be approximated as an exponential distribution with mean values within the range of 200–400 bytes. Fig. 5 illustrates an example where the number of packets belonging to different video configurations is similar.

Therefore, at low traffic loads the resulting network performance will be accurate, while as the traffic load increases the probability of packet drops and the probability of transmission delays will raise exponentially.

4.1.4. Simulation set up

In order to analyze the impact on the QoE of the combined service-level configuration and network-level performance, a number of simulations have been run. Table 3 illustrates the most relevant configuration parameters set up for the case study at service and network levels. For the service configuration, three parameters are considered significant for the user experience. The AMR mode determines the encoding audio quality, so only the three better modes are used. The packetization scheme determines the number of voice frames that are transmitted in a single RTP packet, and has an impact on the header overhead and the latency of the frames. Finally, the size of the dejittering buffer governs the application-level voice frame loss ratio for a same packet jitter pattern. Concerning the UTRAN, the network performance is given by the error pattern at RLC level. The burst pattern is kept constant for all the studies, while the BLER values are varied from good reception conditions – 0.1% of BLER – to severe radio degradations – 30% of BLER. With regard to the core network, the network performance is determined by the traffic intensity and the distribution of the packet size. In this study,

Table 3
Simulation parameters.

AN	
BLER values (%)	0.1, 0.5, 1, 2.5, 5, 10, 20, 30
MBL	1.75
CN	
Traffic intensity (Mbps)	0, 1, 1.5, 1.8, 1.85, 1.875, 1.9, 1.925, 1.95, 1.975
Packet size distribution	Exp(350 bytes)
VoIP	
AMR mode (kbps)	12.20, 10.20, 7.95
Pack. scheme (# frames/packet)	1, 2, 3
Dejittering buffer size (ms)	60, 80, 100, 120

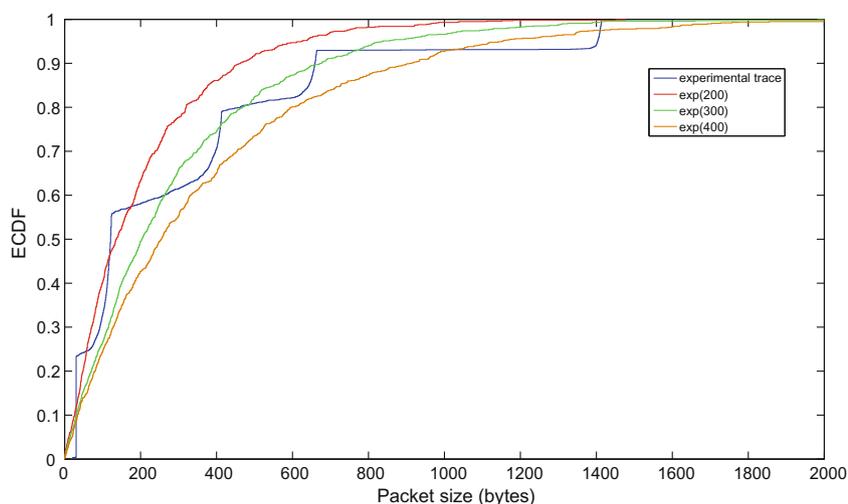


Fig. 5. Exponential approximation of packet size probability distribution for CN.

the packet size pattern is fixed to an exponential distribution with a mean outcome of 350 bytes. Meanwhile, the background traffic load is varied from the unload condition to a heavy loaded network, where the ratio between the arrival and the service capacities are close to one.

In the rest of the document, the VoIP configuration will be represented as a 3-tuple (triple) of parameters (i, j, k) , or as a unique ordered numerical value as given by the following equation:

Service level parameter set = (i, j, k)

Ordered set : x

$$x = (i - 1) \cdot 12 + (j - 1) \cdot 4 + k, \quad (9)$$

$$\begin{cases} i = (1 - 3) & \text{for } \{12.20, 10.20, 7.95\} \text{kbps} \\ j = (1 - 3) & \text{for } \{1, 2, 3\} \text{frames/packets} \\ k = (1 - 4) & \text{for } \{60, 80, 100, 120\} \text{ms of buffer} \end{cases}$$

The service provision scenario has been implemented in OPNET Modeler, including support of a better E-model implementation, two-state Markov modeled bursty radio channel, and the service-level and network-level adaptation capabilities. In order to analyze the impact on QoE of the combined AN/CN effects, all the possible combinations of network states have been configured. In total, a number of 80 combined network states have been tested. For each pair of network parameters, a total of nine possible VoIP configurations are used to transmit the digitalized voice frames, combining all the possible AMR mode and packetization values. These two parameters have a direct impact on the performance results for each considered network state, since they determine the required bitrate and the used packet sizes. Meanwhile, the de-jittering buffer size is considered during the post-processing phase, in order to determine whether VoIP packets arrive on-time to be played out or are finally discarded in the buffer. The simulation time is configured to 300 s for each of the 720 simulations.

The output of the network simulation process is a set of packet traces, each of them including all the VoIP packets successfully received at the user side and the associated e2e delay. Based on this information, the service performance can be inferred by computing the resulting MOS scores. Thus, we finally obtain a total of 36×720 traces with the evolution of the MOS scores over 300 s for each combination of VoIP configuration, AN state and CN state. Based on these results, we analyze the global behaviour of each trace in terms of QoE, in order to determine the optimal service configuration for each network state.

4.2. Selection of best VoIP configurations

The approach described in Fig. 1 allows us to develop an automated decision making process aimed at evaluating the suitability of each VoIP configuration for all the combined AN/CN states. The output of this process is a list that includes the 36 possible VoIP configurations sorted by the obtained MOS_{factor} value, i.e., ordered by the estimated QoE.

Fig. 6 shows the results found where BLER is set up to 5% in the access network, and for all the considered loads in the core network. The y-axis shows the computed MOS_{factor} values, while the x-axis represents all the considered VoIP configurations as given in (9).

As a first result, the QoE-based analysis provides the most suitable VoIP configuration for each combined AN/CN state. These results can be used as a decision making engine intended to drive an automated network-aware dynamic VoIP adaptation procedure.

Additionally, we can easily find how far the different VoIP configurations are in terms of our objective QoE indicator, which can be useful in order to make the final adaptation decision. On one hand, the difference in the MOS_{factor} scale provides an estimation

of how the QoE would be improved if the proposed VoIP adaptation is performed. Thus, it would be another aspect to take into account for evaluating if the proposed switch is worth to be carried out or not. On the other hand, this information could be useful from the standpoint of the optimization of resource usage. For example, if the distance in the MOS_{factor} scale between two possible VoIP configurations is low, it could be preferable the configuration with a lower bitrate requirement – this is, with lower AMR mode or higher packetization scheme.

4.2.1. Service-level decision map

From a similar analysis for all the network states defined, we can find the most suitable VoIP configuration for each considered case. Fig. 7 visually summarizes the results obtained in the whole process.

The values for the VoIP configuration are those given by (9), while the values for the AN states and the CN states are those presented in Table 4.

From these results, we can address the need for considering an adaptation approach which takes into account the combined VoIP configuration capabilities.

- At low BER values of 0.1%, the best configuration for the non-loaded CN is the default AMR-12.20 codec mode with one voice frame per packet, and a de-jittering buffer size of 60 ms, as determined for example in [7]. This minimal buffer is needed mainly due to the transmission effects through the UMTS network, which introduces discrete jitter values.
- The increase of the CN load does not impact the VoIP session until it reaches the 95%. Here, the additional jitter may be mitigated decreasing the required bitrate. Since the jitter due to RLC retransmissions is low, the total delay is kept in tolerant values, and the bitrate can be decreased by the packetization scheme instead of downgrading the codec mode. As the CN load keeps increasing, the additional jitter requires an increase of the buffer size for accommodating the packets. In the limit of 98.75% of CN load, the introduced jitter requires a high buffer size of 100 ms. The total delay is thus increased, and as a result the packetization scheme is set up back to the default value.
- For the low-medium BLER values in the range up to 5%, the effects seem similar. We find a range of CN loads where the additional delays are tolerated, so the best scoring VoIP configuration is given by the trade-off between the delay and loss impairments and driven by the value of the buffer size. Then, we find an area where the jitter introduced by the CN is the most relevant impairment, and the transmission bitrate is reduced. We must note that, in this range, the effect of the bitrate reduction is always worse in terms of QoE by downgrading the AMR mode than by increasing the number of voice frames per packet. Finally, there is a third area where the contribution to the total delay is more worthy due to a higher buffer size than due to a higher packetization scheme. As can be observed, the location of these areas is moved towards lower CN states as the BLER value is increased.
- Another relevant finding is that for all these low-medium BLER configurations, the codec mode is barely selected to switch. At 10% of BLER, the AMR-10.20 mode is selected at two different CN loads, while for 20–30% of BLER the best performing codec mode varies between the three considered values of AMR-12.20, AMR-10.20 and AMR-7.95. As well, for these severe radio degradations, the packetization is switched to high values in order to minimize the contribution of the CN to the delay variation. This effect seems to follow the codec adaptation procedure explained in [8].

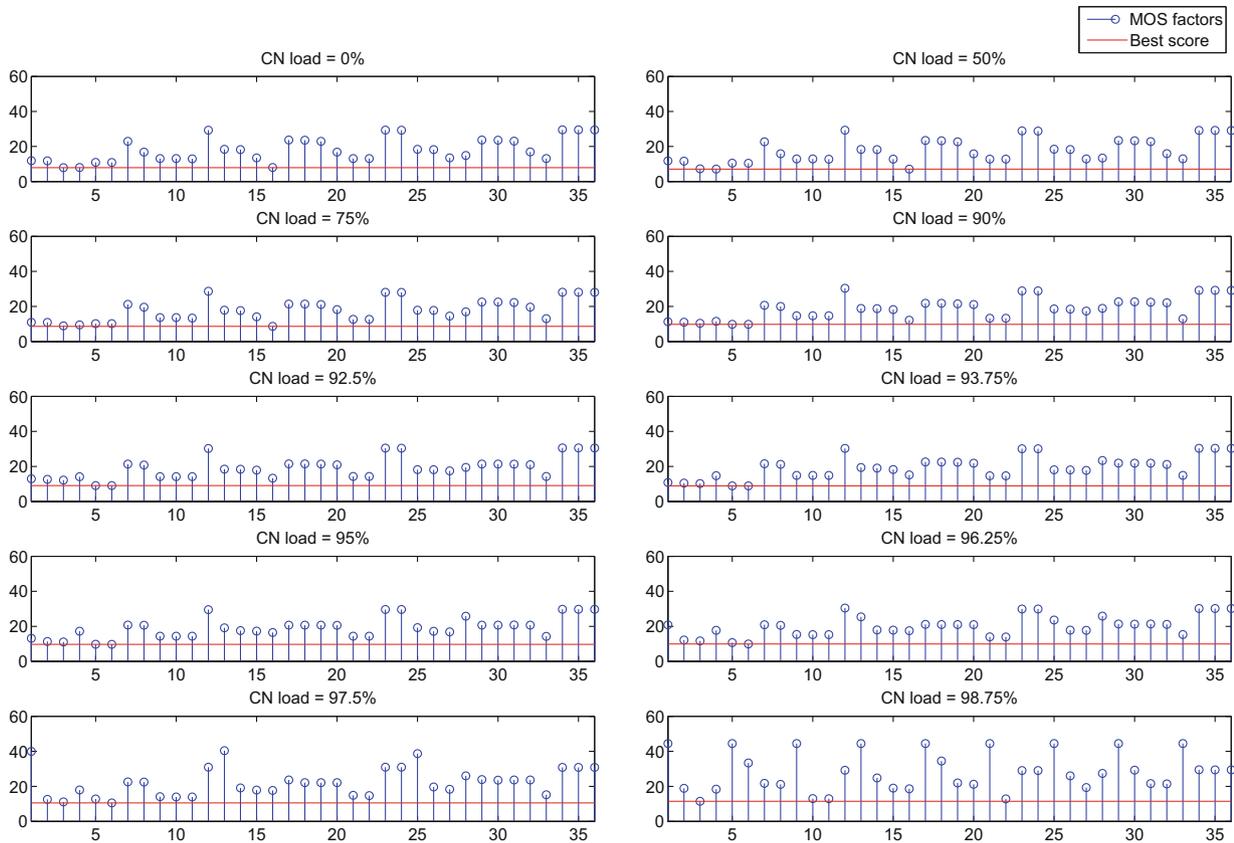


Fig. 6. Evolution of MOS factor for different VoIP configurations vs. best score.

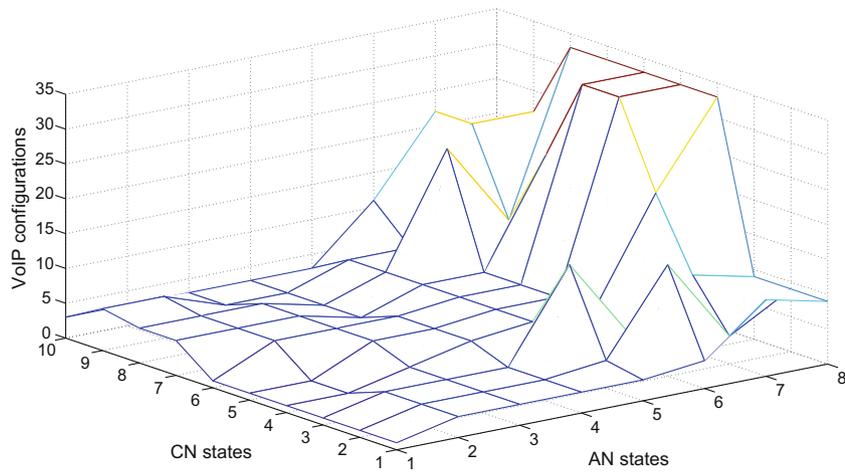


Fig. 7. Map of best VoIP configuration set for different AN/CN states.

Table 4
Best choice of service-level parameters.

AN state BLER (%)	CN state: traffic load (%)									
	0	50	75	90	92.5	93.75	95	92.65	97.5	98.75
0.1	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,2,1)	(1,2,1)	(1,2,2)	(1,1,3)
0.5	(1,1,3)	(1,1,3)	(1,1,1)	(1,1,1)	(1,1,1)	(1,2,1)	(1,2,1)	(1,2,1)	(1,2,2)	(1,1,3)
1	(1,1,3)	(1,1,3)	(1,1,3)	(1,1,3)	(1,2,1)	(1,2,1)	(1,2,1)	(1,2,1)	(1,1,3)	(1,1,3)
2.5	(1,1,3)	(1,1,3)	(1,1,3)	(1,2,1)	(1,2,1)	(1,2,1)	(1,1,3)	(1,1,3)	(1,1,3)	(1,1,3)
5	(1,1,3)	(1,1,4)	(2,1,4)	(1,2,2)	(1,2,2)	(1,2,2)	(1,2,2)	(1,2,2)	(1,2,2)	(1,1,3)
10	(1,1,4)	(2,1,4)	(1,1,4)	(1,2,2)	(1,2,2)	(1,2,2)	(1,2,2)	(2,3,2)	(1,2,2)	(1,3,3)
20	(1,3,3)	(1,1,4)	(1,3,3)	(2,3,1)	(3,3,1)	(3,3,1)	(2,3,1)	(1,2,2)	(2,3,2)	(1,3,2)
30	(1,3,1)	(1,3,1)	(1,3,1)	(3,3,1)	(3,3,1)	(3,3,1)	(3,3,1)	(3,3,1)	(2,3,2)	(1,3,2)

4.3. Analysis of achievable QoE levels

From the previous results, an automated decision system could be able to recommend service-level adaptations based on the actual network performance. In this section, we examine the QoE that the mobile end users could expect as the result of these dynamic adaptations.

Fig. 8 shows the ECDF of the MOS values associated to the VoIP configurations that are selected as the best choices by the automated adaptation system. The figure shows how the CN conditions modify the achievable quality levels for the specific case of a BLER value of 2.5%. In addition, in order to illustrate the impact of both network segment on the achievable quality levels, Fig. 9 represents the results obtained for all the considered AN/CN states in a boxplot. For each AN state, all the traces for the different CN states are shown. From a simple inspection of results, we can observe the high impact of the increasing BLER on the experienced quality.

In order to analyze the actual impact of the CN performance in the QoE, we analyze the amount of samples that are below the three MOS thresholds specified as representative for the user experience. Therefore, Table 5 gathers the maximum and minimum values found in the ECDF of the MOS at those thresholds for all the different BLER values.

As well, Table 6 represents the mean values of the MOS scores found for the selected best VoIP configuration at each combined AN/CN state.

From the data in Tables 5 and 6, we can undertake several relevant conclusions:

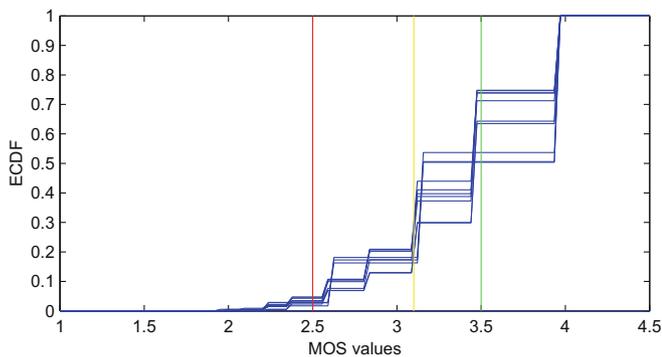


Fig. 8. ECDF of the best VoIP choice for BLER = 2.5% and different CN states.

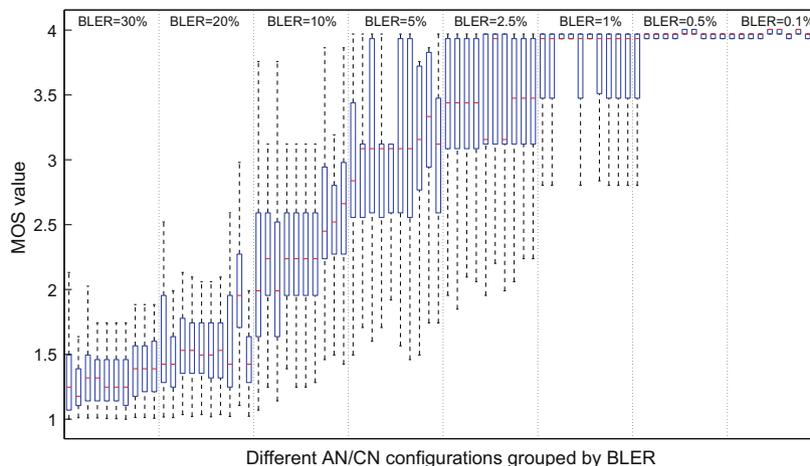


Fig. 9. Boxplot of the achieved MOS for the best VoIP choice per AN/CN state.

Table 5

Percentage of fulfilment at different target QoE thresholds.

BLER (%)	MOS < 3.5		MOS < 3.1		MOS < 2.5	
	Min (%)	Max (%)	Min (%)	Max (%)	Min (%)	Max (%)
0.1	0.71	1.52	0	0.70	0	0
0.5	4.27	12.61	0.67	2.71	0	0
1	9.72	25.83	2.63	6.10	0.37	0.89
2.5	29.85	53.69	12.86	20.92	1.72	4.74
5	63.93	80.96	38.99	55.22	12.27	22.38
10	87.92	97.18	80.57	88.65	42.61	62.71
20	98.42	99.84	98.42	99.34	91.43	94.07
30	99.46	100	99.46	100	98.49	99.50

- For low BLER values, in the range of 0.1–1%, the analysis of the mean MOS values seem to entail accurate QoE levels. The mean MOS is over the 3.5 threshold for all the CN loads considered.
 - Yet, if we focus on the study of the ECDF values in Table 5, we detect that the CN state may be quite relevant for the same experienced BLER. For instance, and increase of the number of samples under the 3.5 threshold from 9.72% to 25.83% may involve a high impact on the QoE.
- For a BLER of 2.5%, the effect of CN is not negligible even based on the simple analysis of the mean MOS values. As can be seen, up to 50% of traffic load the mean MOS is kept above the 3.5 threshold, while for higher loads the mean MOS goes down under that threshold.
 - The ECDF analysis shows a difference of approximately the double in the percentage of samples over and under each MOS threshold between the best and the worse cases.
- Similarly, for 5% and 10% of BLER the CN state determines if the mean MOS values are over or under the thresholds of 3.1 and 2.5, respectively.
- In terms of service acceptability, we must highlight that for BLER of 5% the best case shows a 12.27% of samples under the 2.5 threshold, which is quite a severe degradation for most users. As well, for BLER of 10% and above, the QoE seems to be not acceptable.

Therefore in both cases according to the approach in Fig. 1 we should trigger a network adaptation mechanism in order to jump to other AN/CN state.

5. Lightweight implementation at service endpoints

From previous results, an automated decision making system could infer different adaptation actions in the service provision

Table 6
Mean MOS values for the most optimal VoIP configurations.

AN state BLER (%)	CN state: traffic load (%)									
	0	50	75	90	92.5	93.75	95	92.65	97.5	98.75
0.1	3.9967	3.9967	4.0086	3.9882	4.0073	4.0084	3.9932	3.9932	3.9712	3.9731
0.5	3.8855	3.8855	3.8675	3.8993	3.8941	3.8873	3.8925	3.9000	3.8716	3.8596
1	3.7736	3.7778	3.7426	3.7437	3.7672	3.7244	3.7720	3.7683	3.7272	3.7273
2.5	3.5035	3.5070	3.4241	3.4462	3.4800	3.4716	3.3866	3.3781	3.3797	3.3645
5	3.1079	3.2348	2.9942	2.9974	2.9888	2.9961	3.0018	3.0001	2.9429	2.9164
10	2.6455	2.5640	2.5799	2.3569	2.4056	2.3866	2.3749	2.1917	2.4089	2.2436
20	1.5816	2.0299	1.5838	1.6446	1.6333	1.6321	1.6574	1.6693	1.5469	1.5704
30	1.4284	1.4314	1.4271	1.3384	1.3327	1.3315	1.3417	1.3526	1.2897	1.3290

chain based on the knowledge of the network status. As a result, the problem arises regarding how to take feedback information from the network in order to estimate its performance and make the most optimal decisions.

If the intelligent decision maker is integrated within a network management infrastructure, the AN and CN states could be estimated or monitored from the network elements themselves. However, in an end-to-end service adaptation approach the feedback information must be gathered at the endpoints.

Hereafter, we address the problem of how the different network states can be identified from the user side based on the statistics of the received data. In this case, in addition to the limited information, we must consider the likely limited capacity of a terminal device for processing all this information. Therefore, the identification of network states shall be based on rather simple operations.

From a first inspection of the e2e delay traces, we find the results shown in Fig. 10. Fig. 10a represents the mean values of the e2e delay found for the whole set of traces, while Fig. 10b represents the variance of the e2e delay values of each trace. The x-axis includes all the simulations run for the default values of AMR-12.20 mode and single packetization scheme. The ordering of the traces is per AN state first and per CN state later, this is, the first 10 traces correspond to the 0.1% of BLER for the 10 possible CN loads.

From Fig. 10b, we can state that the variances of e2e delay seems to be a good indicator of the AN state. As can be observed, the found range of values is different and non-overlapping for the different intervals of traces. Thus, in this case we can conclude that the higher variations in the e2e delay are mainly caused by the

effects of the BLER. Each Transport Block erroneously received at the user side is retransmitted by the RNC due to the error recovery function of the RLC AM mode. This retransmission increases the total transmission time of the Transport Block in about 90 ms, which at the same time makes the total packet transmission time increase. Logically, a higher number of retransmissions results on a higher increase on the e2e delay.

The corresponding values of the variance of the delay all over the traces are shown in Table 7. Since they are non-overlapping values, this statistic is proposed for the identification of the AN state.

Once we can identify the BLER value, we focus on the differentiation of the different CN states. As can be seen, the mean e2e delay increases both with the experienced BLER in the AN and with the traffic load in the CN. On one hand, the mean delay increases with the BLER due to a higher number of RLC retransmissions. At higher BLER values, the number of voice frames with at

Table 7
Variance of the delay.

BLER (%)	Range of variance (ms)
0.1	14–52
0.5	103–147
1	204–270
2.5	508–620
5	1050–1230
10	2087–2494
20	4624–5627
30	8654–10341

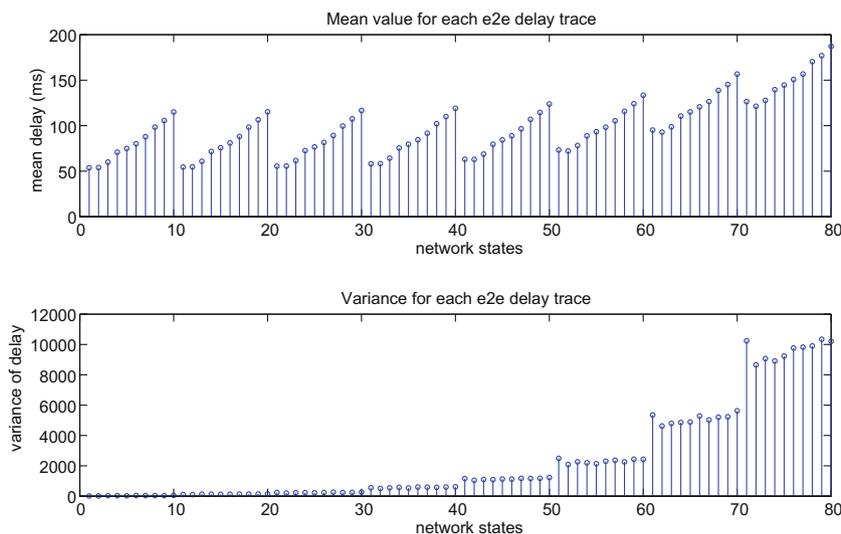


Fig. 10. Mean and variance of e2e delay for different AN/CN states.

least one retransmitted TB increases, contributing to the mean delay with a higher value. In addition, at high BLER values the probability that a TB has to be retransmitted more than once is not negligible. On the other hand, the mean delay increases with the CN load due to the performance of the transmission queues. The associated additional delays are governed by the traffic arrival rate – i.e., the total load over the capacity – and the packet arrival pattern, which in this case is characterized by a packet size that follows an exponential distribution. Since the packet size pattern is fixed in this experiment, the additional contribution of the CN to the mean delay will increase exponentially as the occupancy gets closer to one.

Due to this dual contribution to the e2e delays, these values cannot be used to identify the CN state independently. Therefore, we try to eliminate the contribution of RLC retransmissions by dropping these samples from the e2e delay traces. The resulting truncated trace still shows a great variability of values for the different BLER values within each CN state. The obtained mean values are quite overlapping, and thus are not useful for establishing identification ranges. However, for each individual CN state, the variance of the mean delay values over different AN states is mainly caused by the very high BLER values of 20% and 30%. This effect is produced by the fact that, at very high error ratios, the number of RLC PDU retransmissions becomes very high. Thus, many situations occur where the RLC data to be transmitted to the user side becomes delayed since the TTIs are full with data to be retransmitted, which takes more priority.

To avoid this effect in the identification of AN states at high BLER values, we decide to drop from the delay traces not only the packets with retransmitted segments, but also the following samples. In summary, we find that when two samples per retransmission are dropped the intervals are already non-overlapping. Fig. 11 gathers the evolution of the mean e2e delays for the new truncated delay trace.

In general, the evolution of the computed mean delays seems to follow the trend of the evolution of the traffic load in the CN. One issue that cannot be avoided regards the identification of the two lowest traffic load states. For all the tests performed, the contribution of the CN to the total delay is very similar up to the 50% of traffic load. In most of the cases this effect is negligible for the QoE, but in some specific cases this may result in wrong decisions. For simplicity, those values are also gathered in Table 8.

Table 8

Variation of mean of delays.

CN load (%)	Range of mean delays for the truncated trace (ms)
0	53.6391–54.9374
50	53.8417–55.5828
75	59.8953–60.9208
90	70.7002–71.3646
92.50	74.8481–76.76
93.75	79.9051–80.8707
95	87.0688–87.8273
96.25	97.5119–98.2511
97.50	105.5178–106.0589
98.75	114.3352–114.8968

6. Accuracy of the decision maker

From the results provided in the two previous sections, we are able to propose an automated dynamic VoIP adaptation system which, based on the delay information monitored at the endpoints, is capable of proposing the most suitable VoIP configuration in response to the expected performance of both the AN and the CN. Yet, there is another feature that must be taken into account in such a system, which is the capability of producing real-time responses to the actual network performance levels.

6.1. Responsiveness of the system

The value ranges presented in Tables 7 and 8 are obtained from the analysis of the whole trace, and thus are related to long-term evaluation of the e2e delay values. With the intent to implement a decision maker in an interactive multimedia service, the responsiveness of the system must be analyzed and kept in accurate values.

Thus, in this section we will examine the capability of the proposed network performance identification approach where the measuring interval is reduced. For such a purpose, we evaluate the evolution of the two established statistics over time. Both statistics are continually computed over time when a new delay sample is available, this is, upon the arrival of a new voice packet at the user terminal. Both computations are carried out in a moving-window manner, in order to include the effects of recently received packets and avoid the contribution of distant samples.

Logically, the size of the sample window will have a great impact on the obtained results. Thus, a total of eight window sizes

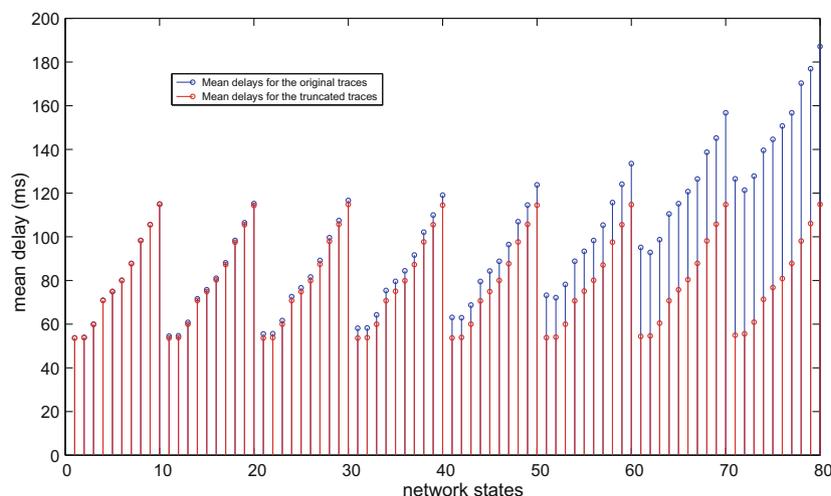


Fig. 11. Mean values for original and truncated delay traces for different AN/CN states.

Table 9

Evolution of window size.

W (number of samples)	50	150	500	750	1000	1250	1500	1750
Considered time (s)	1	5	10	15	20	25	30	35

(W) have been tested, as shown in Table 9. For the case of one voice frame per RTP packet, the relationship between the window size and the considered sample time is governed by the fact that 50 packets per second are generated.

Fig. 13 illustrates the results obtained for the variance of the monitored e2e delays for all the considered window sizes. Each subplot includes the results for all the traces, grouped in colours in a per AN state basis. As can be observed, the different ranges appear more clearly as the window size is increased.

Likewise, Fig. 12 shows the results of similar experiments for the mean values of the truncated delay traces. In this case, the different traces are grouped in different colours based on their associated CN load. The simple visual inspection shows that the different ranges seem non-overlapping even for low window sizes.

As a result, the relationship between the window size and the width of the ranges can be seen as a trade-off. Translated to the QoE plane, this trade-off determines the relation between the responsiveness of the system – this is, how quick the adaptation actions can be produced – and the accuracy of the decisions – this is, the suitability of the network identifications and the consequent VoIP adaptation proposals.

6.2. Accuracy of network state estimations

Based on the evolution of the computed values for the selected statistics, we develop a module that continuously provides estimations of the network states over the time. As a compromise of

accuracy and responsiveness, we show the results found for the window size of 250 samples.

From Fig. 10b, it seems that the CN state can be suitably followed with this approach. However, the evolution shown in Fig. 10b for the variance does not allow us to foresee a good identification outcome. From a detailed analysis of the variance over time, we detect that for $W = 250$ the percentages of successful identifications are between 13.65% and 37.19%. The percentage of wrong decisions is always under the 1% for BLER values of 10–30%, and this percentage increases as the BLER comes lower. Yet, except for the lowest BLER values under the 1%, the percentage of good identifications is always superior to the percentage of wrong identifications.

Therefore, we establish the window size of 250 samples as a reasonable solution for the accuracy vs. responsiveness in the system. This window value is not only a relevant factor in the trade-off between responsiveness and accuracy, but it has also determines the overload that the endpoint has to support for inferring network states. Since we find good results for 250 samples, the number of delay values to be stored and the related number of operations is very limited, even for mobile devices.

For low BLER network states, up to 1%, this approach is not able to provide reliable estimations. Yet, these are best network performance conditions, and the final impact on the QoE may be limited. This effect is later analyzed, where we study the impact of possible wrong decisions on the user experienced quality.

In summary, the logic for the estimation of the network states is implemented as follows:

- For each new delay sample, the two selected statistics are updated. Only the past samples within the window size are considered for the computation.
- If any of the obtained values fall into one of the established ranges, the corresponding network state is updated. Otherwise, both network states are kept as in the previous iteration.

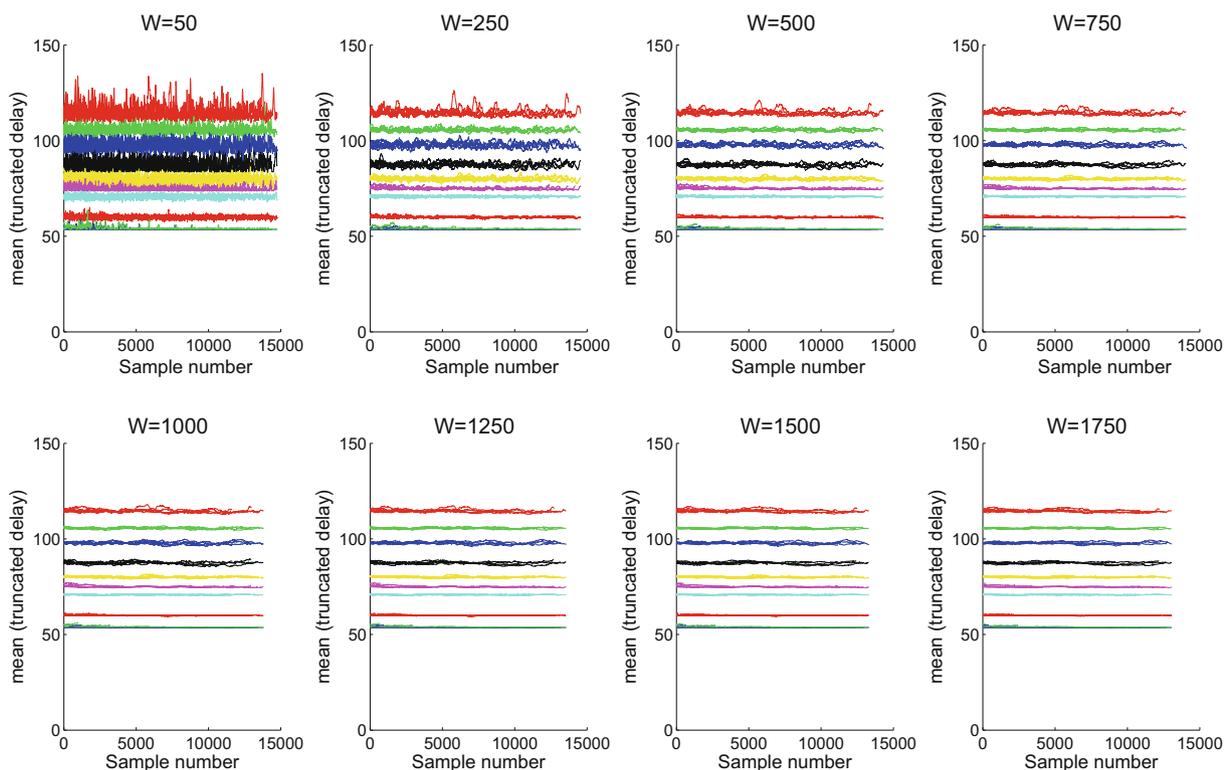


Fig. 12. Computed e2e delay for truncated traces over time for different window sizes.

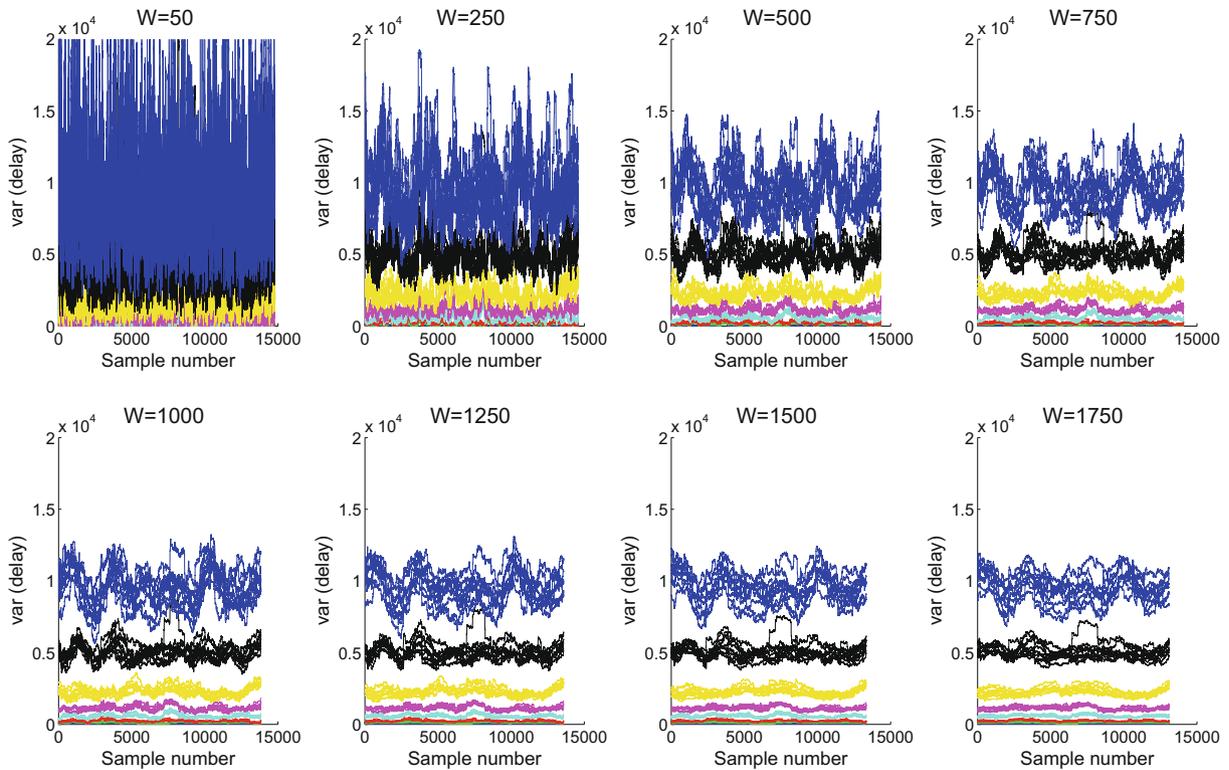


Fig. 13. Computed variance of the e2e delay for different window sizes.

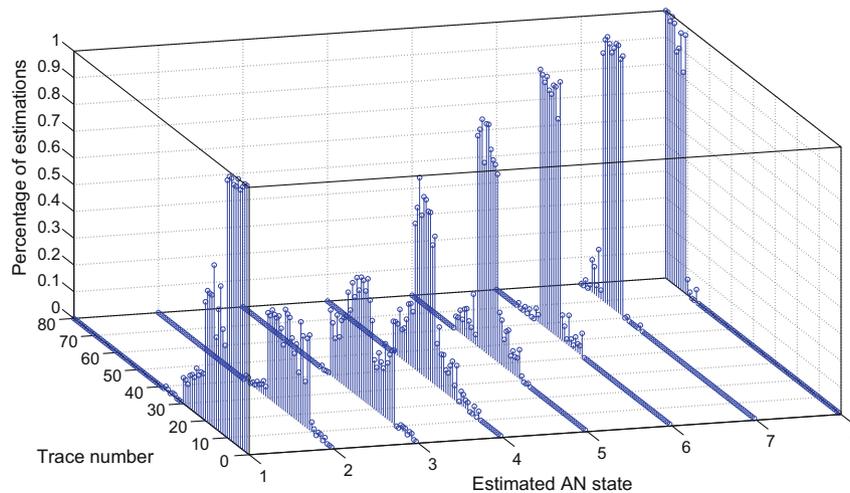


Fig. 14. Ratio of AN state proper estimations.

Fig. 14 illustrates the ratio of good and wrong decisions in the estimation of the AN state. The delay traces are ordered per AN state first, and then per CN state. As expected, the worst results are found for BLER values of 0.5% and 1%, which are those traces from 11 to 30. Even in those worse cases, we can see that the wrong decisions mainly fall in the adjacent BLER states, which may limit the grade of failure in the decision making. For high BLER values, the percentage of good estimations is very high, which allows a quick identification of severe network degradations in order to perform quick adaptations.

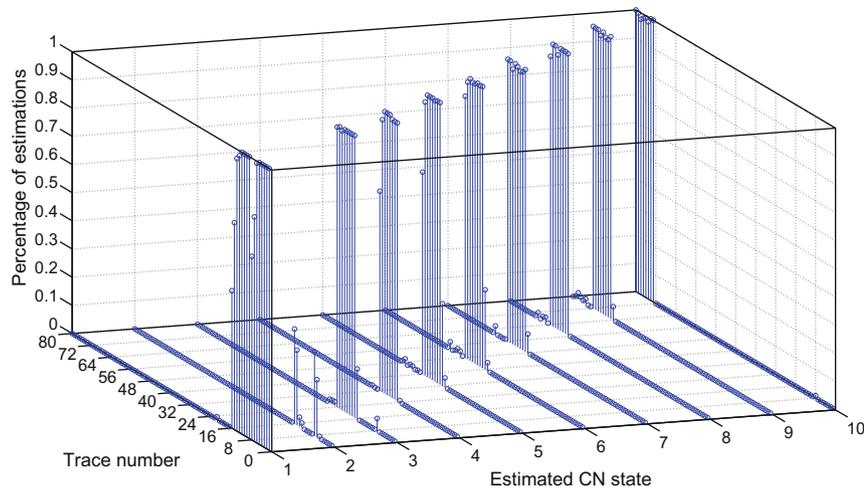
Table 10 gathers the Root Mean Square Error (RMSE) values associated to the estimations. For each combination of

AN and CN states, the table shows the RMSE of the BLER estimation.

Similarly, Fig. 15 represents the ratio of CN state estimations for all the simulations. In this case, the delay traces have been sorted per their configured CN state first, and based on the AN state secondly. Again, the expected behaviour is followed. For the first two CN states – trace numbers from 1 to 16 – the contribution to the total delay cannot be differentiated and both states are identified as the non-loaded network state. This effect can be clearly observed in Table 11, where the RMSE for the second CN state tends to one for all the cases. This means that most of the times the CN load state of 50% is estimated as a non-loaded CN state. The rest of the CN loads are quite well estimated.

Table 10
RMSE of estimated AN states.

AN state BLER (%)	CN state: traffic load (%)									
	0	50	75	90	92.5	93.75	95	92.65	97.5	98.75
0.1	0	0	0.2619	0.2931	0	0.2928	0.2928	0.2623	0.1319	0.1917
0.5	0.8528	0.8142	0.965	0.8163	0.7875	0.9498	0.8992	0.9454	0.9274	0.8107
1	1.0884	1.0545	0.992	1.0465	0.9702	0.9728	0.9976	0.9023	0.9605	0.9647
2.5	0.7441	0.8017	0.6837	0.7177	0.6977	0.6023	0.7144	0.6255	0.6664	0.7183
5	0.5585	0.5275	0.508	0.5068	0.3707	0.3768	0.5415	0.3706	0.4261	0.4598
10	0.1959	0.4277	0.263	0.2644	0.3313	0.3186	0.2379	0.2908	0.2511	0.2215
20	0.1649	0.2117	0	0	0.1508	0.2148	0.1385	0	0.1309	0.2705
30	0.184	0.4235	0.1968	0.32	0.3493	0.1481	0.0603	0.1647	0	0

**Fig. 15.** Ratio of CN state proper estimations.**Table 11**
RMSE of estimated CN states.

AN state BLER (%)	CN state: traffic load (%)									
	0	50	75	90	92.5	93.75	95	92.65	97.5	98.75
0.1	0	1	0	0.1174	0	0.0388	0	0.0438	0	0
0.5	0	1	0	0.1177	0	0.0585	0.1114	0	0.1253	0
1	0	1	0	0.1095	0.0603	0	0.1269	0.0117	0.1144	0.1272
2.5	0	1	0	0	0	0.0884	0.0722	0	0.0083	0.0803
5	0	0.9912	0	0.0143	0	0.0892	0	0.166	0.1282	0.1182
10	0.0848	0.9835	0.1924	0	0	0	0.1564	0	0.0341	0.1511
20	0.4566	0.857	0.0166	0.1834	0.1713	0.1319	0	0	0.0341	0.0763
30	1	0.8153	0.0811	0.5427	0.5331	0.2691	0	0.2187	0.0117	0.0083

6.3. Accuracy of VoIP-level decision making

Finally, we test the behaviour of the decision maker itself. As previously introduced, the goodness of the estimation of the network performance is a mean for the final objective of the system, which is not other than keeping the QoE levels in the optimal values for the different network conditions.

Therefore, the percentage of good or wrong estimations of network states is only an indicator of the system performance. As well, the distance between a wrong estimation and the actual value must be considered, this is, how erroneous the estimation is. It is expectable that the performance of one specific VoIP configuration follows a continuous trend among adjacent network states.

First, we obtain the percentage of time that the decision maker proposes as the best choice the VoIP configuration that corresponds for each configured combined AN/CN state. In this case, the RMSE is not a valid indicator for understanding how wrong

the decisions are. The sorting of the 36 possible VoIP configurations considered does not follow a logical impact on QoE, and does not even include a sense of proximity between VoIP configurations. Because of that, Table 12 gathers just the percentage of times when the optimal VoIP configuration is selected by the system.

Due to the existence of flat areas in the decision map (see Fig. 7), the percentage of correct decisions is increased over the expected value if we take into account the ratio of accurate combined AN/CN state identifications.

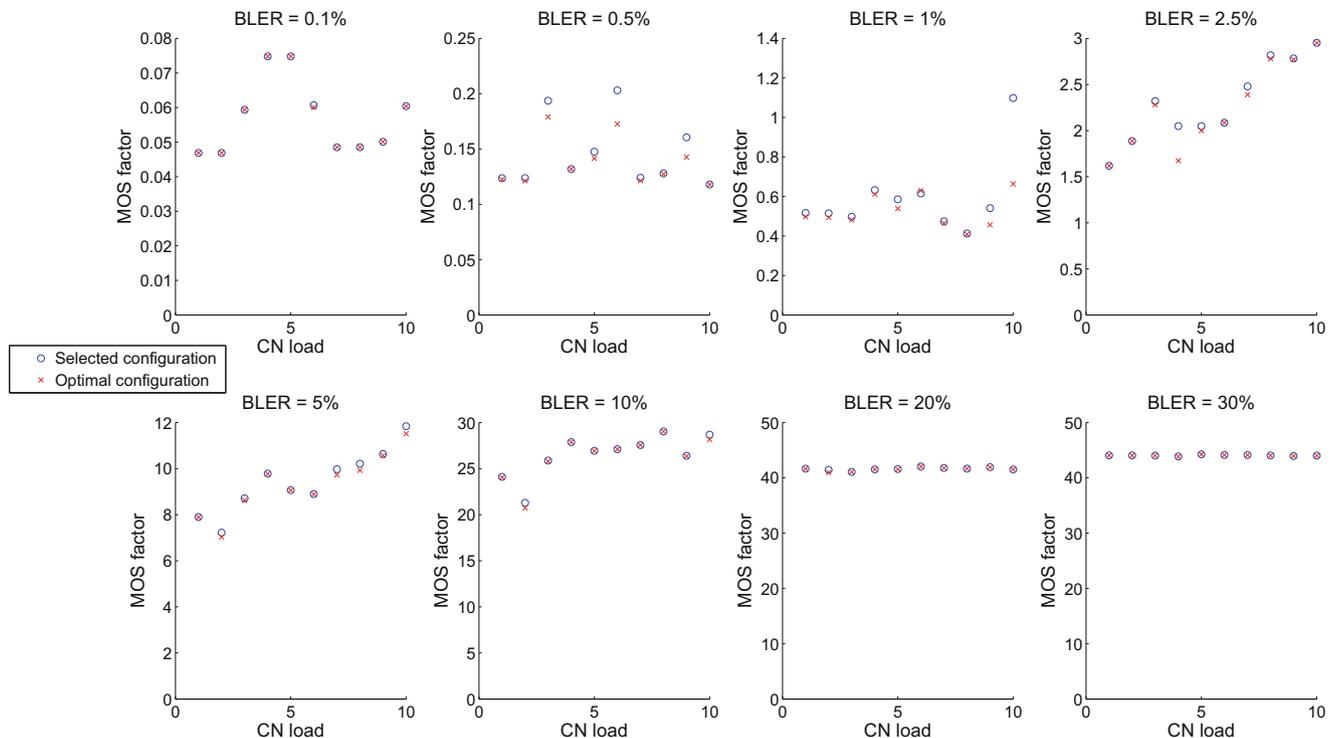
Yet, there are several combined states where the ratio of correct decisions is not acceptable. For example, when the BLER is set to 1% and the CN load fluctuates between 90% and 92.5%, the proposed VoIP configuration is the optimal one only for half the time. Even more, for CN loads of 50% and BLER values between 5% and 10%, the optimal VoIP configuration is rarely selected.

Thus, we must study what happens for the rest of the time, where the VoIP configuration proposed by the decision maker is

Table 12

Percentage of correct decisions for the VoIP configuration.

AN state BLER (%)	CN state: traffic load (%)									
	0	50	75	90	92.5	93.75	95	92.65	97.5	98.75
0.1	100	100	98.28	98.28	100	96.57	100	100	99.99	100
0.5	66.11	60.71	45.17	78.89	72.37	39.25	95.67	95.46	77.09	100
1	83.02	83.10	84.92	41.61	53.17	86.61	83.94	86.75	59.09	98.40
2.5	100	99.93	93.02	62.88	90.23	84.99	59.32	72.95	81.30	100
5	91.52	3.45	74.19	82.75	87.77	88.93	82.07	86.27	87.90	90.42
10	95.44	2.20	93.09	100	100	100	100	91.54	98.75	92.81
20	77.10	26.37	100	96.64	97.70	100	98.08	100	99.88	92.68
30	96.61	82.07	95.47	91.12	100	99.96	99.64	92.55	99.99	99.99

**Fig. 16.** Comparison of the achieved QoE levels: actual decisions vs. optimal configuration.

under the optimal solution. For such a purpose, we have to analyze the distance in the QoE scale between each selected VoIP configuration and the optimal one. And that distance is included in the concept of MOS factor introduced in this paper.

Fig. 16 illustrates the results on the QoE dimension. For all the possible combinations of AN and CN states defined, the figure shows the comparative results of the optimal solution and the actual selection. The optimal solution represents the VoIP configuration with the minimum MOS factor among all the 36 possibilities. For the actual selection, the presented MOS factor is the mean of the MOS factors for all the selected VoIP configurations, weighted by the percentage of time that each VoIP configuration is selected by the decision maker.

For example, BLER = 5% and $CN_{load} = 50\%$ is most of the time estimated as BLER = 5% and $CN_{load} = 0\%$, which is another VoIP configuration that behaves very similar at MOS scale.

7. Conclusions

The results shown in this paper constitute an advance towards a comprehensive QoE-driven cross-layer management system for VoIP over 3G UMTS services. Two features are critical for achieving

the best possible QoE levels in a dynamic self-configuring approach: the understanding of the different possible sources of degradation and the inclusion of the network-awareness in the decision making for the selection of the most optimal configuration of the service provision chain.

The first contribution is related to the specification of a network-based service-level adaptation mechanism, aimed at maximizing the QoE levels for varying network performance states. In Section 4.2 we show how the most optimal VoIP configuration evolves for different combined AN/CN states. Thus, from the knowledge of the current performance levels related to the different network segments, the best performing VoIP configuration is selected. The analysis of the decision map allows us to address the need for taking into account different service-level adaptations at the same time, even at both endpoints of the communication.

The specific knowledge of the AN and CN network states allows us to infer the combined impact into the QoE, and the combined service adaptation endows the system with the capability to propose the reaction that better mitigates the degradation effects into the QoE scale. The use of this kind of decision maps allows an easy implementation of the logic. Although not all the network states are covered, the evolution of the selected VoIP configurations shows the tendencies in the optimal QoE curve. The study of the

service acceptability in Section 4.3 determines if the proposed dynamic service-level adaptations are worth enough, in relation to the QoE requirements established for the user. Thus, from the analysis of acceptability the decision making process is able to infer if network adaptations are required. In that case, based on the knowledge of the available network-level adaptation procedures, the system is able to detect the required network state and the associated optimal VoIP configuration.

With regard to the actual performance of the intelligent adaptation system, we show how the logic can be implemented in a lightweight mode. In Section 5 we show how the different network states can be inferred from a simple analysis of the packet delays. Thus, this feedback information allows the system to make the adaptation decisions based on the information available in the simple decision map. This feature endows the required edge-based intelligence to the system.

The AN state is inferred from the variance of the e2e delays, and this the inter-arrival packet can be used as well. Yet, the CN state is identified based on absolute delay values, so at least a first estimation of the network delay is required. Concerning the responsiveness issues, the estimation of network states from the endpoints entails a trade-off between the accuracy of the estimations and the quickness in the reactions. Section 6 shows how a quite reliable decision system can be achieved with a time interval of 5 s. Compared to the state-of-the-art implementations, this approach behaves quite well for the management of interactive communications.

The worst performance of the system is found for low BLER values, since the required time for capturing changes in the RLC loss pattern is higher. This characteristic is inherent to the nature of the UMTS performance, and can only be overtaken with a wider monitoring interval. Yet, with a wider sample window the system would lose the capability to provide fast reactions to severe degradations. Since the network identification procedure is kept quite simple, the implementation of two simultaneous computation processes would be feasible, one for the long-term and fine identification of slight changes, and the other one for the short-term and quick detection of severe impairments.

Acknowledgment

The research leading to these results has received funding from the European Community's Seventh Framework Programme FP7/2007–2013 under Grant Agreement No. 214751//ICT-ADAMANTI-UM/.

References

- [1] S. Maniatis, E. Nikolouzou, I. Venieris, et al., QoS issues in the converged 3G wireless and wired networks, *IEEE Communications Magazine* 40 (8) (2002) 44–53.

- [2] E. Myakotnykh, R. Thompson, Adaptive speech quality management in voice-over-IP communications, in: *Fifth Advanced International Conference on Telecommunications (AICT 2009)*, vol. 0, Venice/Mestre, Italy, May 24–28, 2009, pp. 64–71.
- [3] Third Generation Partnership Project, TS26.975; Performance Characterization of the Adaptive Multi-Rate (AMR) Speech Codec, 2008. Available from: <<http://www.3gpp.org/ftp/Specs/html-info/26975.htm/>>.
- [4] F. Poppe, D. De Vleeschauwer, G. Petit, Choosing the UMTS air interface parameters, the voice packet size and the dejittering delay for a voice-over-IP call between a UMTS and a PSTN party, in: *Twentieth Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE INFOCOM 2001)*, vol. 2, 2001, pp. 805–814.
- [5] I. Mkwawa, E. Jammeh, L. Sun, A. Khan, E. Ifeachor, Open IMS core with VoIP quality adaptation, in: *Fifth International Conference on Autonomic and Autonomous Systems (ICAS 2009)*, vol. 0, IEEE Computer Society, Valencia, Spain, April 20–25, 2009, pp. 295–300.
- [6] I. Curcio, J. Kalliokulju, M. Lundan, AMR mode selection enhancement in 3G networks, *Multimedia Tools and Applications* 28 (3) (2006) 259–281.
- [7] A. Barbaresi, A. Mantovani, Performance evaluation of quality of VoIP service over UMTS-UTRAN R99, in: *IEEE International Conference on Communications (ICC 2007)*, vol. 7, Glasgow, Scotland, 2007, pp. 634–639.
- [8] T. Hoßfeld, A. Binzenhöfer, Analysis of Skype VoIP traffic in UMTS: end-to-end QoS and QoE measurements, *Computer Networks* 52 (3) (2008) 650–666.
- [9] S. Khan, S. Thakolsri, E. Steinbach, W. Kellerer, Qoe-based cross-layer optimization for multiuser wireless systems, in: *18th ITC Specialist Seminar on Quality of Experience*, Karlskrona, Sweden, 2008, pp. 63–72.
- [10] S. Thakolsri, S. Khan, E. Steinbach, W. Kellerer, Qoe-driven cross-layer optimization for high speed downlink packet access, *Journal of Communications, Special Issue on Multimedia Communications, Networking and Applications* 4 (9) (2009) 669–680.
- [11] ITU-T, Recommendation G.107; The E-model, a computational model for use in transmission planning (2003).
- [12] J. Fajardo, F. Liberal, N. Bilbao, Study of the impact of UMTS best effort parameters on QoE of VoIP services, in: *Fifth International Conference on Autonomic and Autonomous Systems (ICAS 2009)*, vol. 0, IEEE Computer Society, Valencia, Spain, April 20–25, 2009, pp. 142–147.
- [13] ITU-T, G.113 Appendix I (05/2002); Transmission impairments due to speech processing. Provisional planning values for the equipment impairment factor le and packet-loss robustness factor Bpl (2002).
- [14] L. Sun, E. Ifeachor, Voice quality prediction models and their application in VoIP networks, *IEEE Transactions on Multimedia* 8 (4) (2006) 809–820, doi:10.1109/TMM.2006.876279.
- [15] A. Raake, Short- and long-term packet loss behavior: towards speech quality prediction for arbitrary loss distributions, *IEEE Transactions on Audio, Speech & Language Processing* 14 (6) (2006) 1957–1968.
- [16] Third Generation Partnership Project, TR25.993; Packet switched conversational multimedia applications; Default codecs (2008). Available from: <<http://www.3gpp.org/ftp/Specs/html-info/26235.htm/>>.
- [17] A.G.O. do Nascimento, E. de Sousa Mota, S.J.B. de Queiroz, E. do Nascimento Jr., An alternative approach for header compression over wireless mesh networks, in: *International Conference on Advanced Information Networking and Applications Workshops (AINA 2009)*, vol. 0, IEEE Computer Society, Bradford, UK, 2009, pp. 165–169, doi:10.1109/WAINA.2009.141. Available from: <http://ieeecomputersociety.org>.
- [18] Third Generation Partnership Project, TR25.993; Typical examples of Radio Access Bearers (RABs) and Radio Bearers (RBs) supported by Universal Terrestrial Radio Access (UTRA) (2005). Available from: <<http://www.3gpp.org/ftp/Specs/html-info/25993.htm/>>.
- [19] W. Karner, O. Nemethova, P. Svoboda, M. Rupp, Link error analysis and modeling for video streaming cross-layer design in mobile communication networks, *ETRI Journal* 29 (5) (2007) 569.
- [20] J. Fajardo, F. Liberal, N. Bilbao, Impact of the video slice size on the visual quality for h.264 over 3g umts services, in: *Sixth International Conference on Broadband Communications, Networks, and Systems (BROADNETS 2009)*, Madrid, Spain, 2009, pp. 1–8.