

A Study of Real-Time Packet Video Quality using Random Neural Networks

Samir Mohamed, *Student Member, IEEE*, Gerardo Rubino

Abstract

An important and unsolved problem today is the automatic quantification of the quality of video flows transmitted over packet networks. In particular, the ability to perform this task in real time (typically for streams sent themselves in real time) is specially interesting. The problem is still unsolved because there are many parameters affecting video quality and their combined effect is not well identified and understood. Among these parameters, we have the source bit rate, the encoded frame type, the frame rate at the source, the packet loss rate in the network, etc. Only subjective evaluations give good results but, by definition, they are not automatic. We previously explored the possibility of using Artificial Neural Networks to automatically quantify the quality of video flows and we showed that they can give results well correlated with human perception.

In this paper, our goal is twofold: First, we report on a significant enhancement of our method by means of a new neural approach, the Random Neural Network model, and its learning algorithm, which both offer many advantages. Second, we follow our approach to study and analyze the behavior of video quality for wide range variations of a set of selected parameters. This may help in developing control mechanisms in order to deliver the best possible video quality given the current network situation, and in better understanding of QoS aspects in multimedia engineering.

Keywords

Packet video, Random Neural Networks, Real-time video transmission, Video quality assessment, Video signal characterization.

I. INTRODUCTION

EVEN though the main concern of multimedia QoS is to maximize the quality of the data delivered to the destination point, most of the current discussion is concentrated on finding a way to keep within certain limits some of the network parameters (e.g., packet loss rate and delay variation), and little attention has been paid to the quality perceived by the end-users of the applications running over the network.

There are several parameters that affect the quality (quality-affecting parameters) of video transmission over packet networks. We can classify them as follows:

- Coding and compression parameters: They control the amount of quality losses that takes place during the encoding process; so they depend on the type of the encoding algorithm (MPEG, H26x, etc.), the output bit rate, the frame rate (the number of frames per sec.), the temporal relation among frame types, etc.;
- Network parameters: They are a result of packetization of the video stream and the transmission through the network, such as the packet loss rate, the loss distribution, the delay, the delay variation (jitter), etc.;

Manuscript received _____

IRISA/INRIA, Campus de Beaulieu, 35042 Rennes Cedex, France. E-mails: Samir.Mohamed@irisa.fr, Gerardo.Rubino@irisa.fr

This work was partially supported by the European ITEA Project 99011 "RTIPA" (Real-Time Internet Platform Architectures).

- Other parameters like the nature of the scene (e.g. amount of motion, color, contrast, image size, ...) can also have an impact on the human perception of the flow.

The analysis of the video quality can be done using either *objective* tests or *subjective* ones. Objective tests are always explicit functions of measurable parameters related to the encoder or to the network [1]. Subjective tests are based on evaluations made by human subjects under well defined and controlled conditions [2], [3]. Obviously, the reference is the end-user's perception, which is directly captured by subjective tests. Concerning available objective tests, it is well known that they do not always correlate very well with human perception [4], [5], [6].

Previous studies either concentrated on the effect of network parameters without paying attention to encoding parameters, or the contrary. Papers that consider network parameters and use subjective tests for the evaluation restrict the study to only one or two of the most important ones. For example, the Loss Rate and Consecutive Lost Packet metrics are studied in [7], while [8] studies mainly the effect of Bit Rate and [9] works only on the effect of Frame Rate, *etc.* Other examples are [10], [5], [6], [11] (see Section II for more details). The main reason for this is the fact that subjective quality tests are expensive to carry out. Moreover, to analyze the combined effect, for instance, of three or four parameters, we need to build a very large set of human evaluations in order to reach a minimal precision level. Concerning objective approaches, there is no previously published objective quality test that can take into account the direct influence on the quality of the whole set of parameters simultaneously.

Some existing objective methods are MSE (Mean Square Error) or PSNR (Peak Signal to Noise Ratio) which measure the quality by simple pixel-to-pixel comparisons. There are other more complicated methods such as the moving picture quality metric (MPQM) and the normalized video fidelity metric (NVFM) [8], [12]. A state of the art for objective video quality assessment methods is [1]. Subjective quality assessment methods [13], [5], [2], [14], [3], measure the overall perceived video quality. They are carried out by human subjects. The most commonly used one for video quality evaluation is the Mean Opinion Score (MOS) [2], [3], recommended by the ITU. It consists of having a set of subjects viewing the distorted video sequences in order to rate their quality, according to a predefined quality scale. That is, human subjects are trained to *build* a mapping between a set of processed video sequences and the quality scale.

Although MOS studies have served as the basis for analyzing many aspects of signal processing, they present several limitations: a) very stringent environments are required; b) the process can not be automated; c) it is very costly and time consuming, making very difficult to repeat it frequently. Consequently, it is impossible to use it in real-time quality assessment. On the other hand, the disadvantages of objective methods are: a) they do not always correlate well with human visual perception¹; b) they require high calculation power, and are time consuming (they usually operate at the pixel level); c) it is very hard to adapt them to real-time quality assessment, as they work on both the original video sequence and the transmitted/distorted one; d) as stated before, it is difficult to build a model that take into account the effect of many quality-affecting parameters, specially network parameters.

¹By the way, this claim comes from comparing the results to those obtained from subjective methods.

In this paper, we propose an approach that, first of all, can evaluate the combined effect of an arbitrary number of parameters on the quality of a video sequence. Moreover, our approach has two supplementary properties: (i) it correlates well with the results obtained from subjective tests (because, as we will see, it is in fact based on them) and (ii) it can work automatically and in real time. Specifically, we build a tool that takes as input the values of a set of parameters associated with the encoder and with the network used to transmit video stream, and correspondingly quantifies the video quality. The tool is based on a neural network trained with the results of previously performed subjective tests, in which wide ranges of the selected parameters and real network conditions are considered. In addition, we use this tool to study and analyze the impact of certain quality-affecting parameters (the source bit rate, the encoded frame type, the frame rate, the packet loss rate, and the loss burst size) on real-time video quality.

In [15], a method to evaluate video quality as a function of certain quality-affecting parameters method based on standard Artificial Neural Networks (ANN) was proposed. Here, we report on a significant enhancement of this method, obtained by using a different type of neural network called Random Neural Network (RNN). This recently invented tool [16], [17], [18] appears to capture with higher accuracy and in a more robust way the function mapping the various involved parameters to the quality metric. This is partly due to nice mathematical properties exhibited by RNN, which makes one of their main differences with ANN.

The organization of the rest of our paper is as follows: in Section II, we give a brief overview of related works. Section V summarizes our previously proposed model to evaluate real-time video quality, introduces the theory of RNN, compares both RNN and ANN, and mentions some benefits of our model. Section VI, presents the subjective quality tests. In Section VII, we describe the quality-affecting parameters and the MOS test we carried out. Finally, in Section VIII, we study the impact of the quality-affecting parameters on video quality.

II. RELATED WORKS

In [10], the authors study the effect of both loss and jitter on the perceptual quality of video. They argue that, if there is no mechanism to mask the effect of jitter, the perceived quality degrades in the same way as it degrades with losses. In [19], [7], [5], [20], [14] and [6], the effect of audio synchronization on the perceived video quality is analyzed (for instance, by quantifying the benefits of audio synchronization on the overall quality of the flow). The main goal of [9] is to study the effect of the frame rate for different standard video sequences on the overall perceived quality. A related work is [11] where the effect of FEC and Frame Rate on the quality is the subject of the study. The work presented in [21] is a study of the packet loss effects on MPEG video streams. The authors consider the effect of loss rate on the different types of MPEG frames. In [8] and [1], a study of the effect of bit rate on the objective quality metrics (PSNR, NVFM, and MPQM) is presented. The effect of the number of consecutively lost packets on video quality is analyzed in [7]. The authors of [22] study the effect of packet size and the distribution of I-frames in the layered video transmission over IP networks. In [13] the analysis goes deeper: the authors present a study of the effect of motion on the perceived

video quality. In [23], a dynamic bit-rate prediction method is proposed. It consists of predicting the bit-rate for future frames based on the past information and dynamically change the quantization parameter if the estimated bit rate exceed certain threshold so that the encoder output's objective quality remains constant. In [24], another objective video quality assessment method is proposed, based on the evaluation of the quality of each frame separately. The idea is to compute a weighted mean of these individual frame values, taking into account the short-term characteristics of human memory. The work in [25] presents a methodology for video quality assessment using objective parameters based on image segmentation. An image encoded by MPEG-2 is segmented into three regions: plane, edges, and texture; then, a set of objective parameters is assigned to each region. After that, a perceptual-based model is defined by computing the relationship between objective measures and results of subjective tests. In [26], a de-jittering scheme is proposed to compensate from the effect of delay variation for the transport of MPEG-4/2 video streams. Furthermore, in [26], an algorithm to protect the packet-loss resilience is proposed. It is based on the switching between Inter/Intra macro-bloc encoding. Although the main goal of the works presented in [4], [5] and [6] is to develop objective method that give good correlation with subjective method, their obtained results clearly shows that globally the two types of outputs do not always correlate well.

In previous work, we showed how to use ANN to measure in real time audio quality when this audio is transmitted over a packet network [27]. Based on this technique, we developed a new control mechanism that permits a better use of the given bandwidth and the delivery of the best possible audio quality given the current network situation [28]. Then, we proposed a tool to evaluate video quality in real time when this video is subjected to wide range variation of certain network and video parameters [15].

Concerning the specific neural approach followed in the present paper, based on the RNN model, it is important to mention that RNN are used in several related problems. For example, they are used in video compression with compression ratio that goes from 500:1 to 1000:1 for moving gray-scale images and full-color video sequences respectively [29]. They are also used as decoders for error correcting codes in noisy communication channels, which reduces the error probability to zero [30]. Furthermore, they are used in a variety of image processing techniques which go from image enlargement and fusion to image segmentation [31]. A survey of RNN applications is given in [32].

III. SUBJECTIVE QUALITY TESTS

Given that we have P parameters to be taken into account, namely ($\mathcal{P} = \{\pi_1, \pi_2, \dots, \pi_P\}$). For each parameter π_p there should be n_p different values. If we choose the whole set of configurations of the P parameters, the total number of samples may be very large ($n_1 \times n_2 \times \dots \times n_P$). To do subjective quality test for this number it may be impractical. Fortunately, the NN does not need all these samples to learn the problem, it is its goal to predict the quality in the missed configurations. Therefore, we propose the following method to generate the minimum number of samples needed to train and test the neural networks. We should define default values for all the parameters, $\mathcal{P}_l = \{\pi_{01}, \pi_{02}, \dots, \pi_{0P}\}$.

These values should be the most frequent observed value. We change the values of two parameters at a time and give the default values to the others. We repeat this step for all the parameters. In this way, a list of configurations was defined, containing many duplications. Once these duplications are removed, the remaining constitute the minimum samples required to train and test the NN.

The resulting database is divided into two parts as previously mentioned. The length of the testing database may be about 10% of the total samples (should be selected randomly from the whole database). The remaining samples represent the training database.

To evaluate the quality of speech, audio or video systems (codec, telecommunication, television pictures, IP telephony, etc.), a subjective quality test is used. In this test, a group of human subjects is invited to judge the quality of the sequence under a predefined system conditions (distortions). There are several recommendations that specify strict conditions to be followed in order to carry out the subjective test. These recommendations are specific to the media type (speech, audio, video, or multimedia in general) in study.

In this Section, we focus on how to use subjective quality tests in our method. In addition, we resume the recommendations of ITU regarding this point.

Depending on the media type (speech only, audio in general, video only, speech and video, speech and audio, multimedia in general [speech, audio and video]), there exists some ITU recommendations that define how to carry out subjective quality tests.

ITU-T Rec. P.800 [?] presents some methods for subjective determination of transmission speech quality. Another recommendations that deals with speech are . . .

For audio in general, ITU-R Rec. BS.1116 [?] recommends methods for the subjective assessment of audio systems. ITU-R Rec. BS1284 [?] presents the general requirements for subjective assessment of sound quality. How to assess subjective audio quality mixed with pictures is the subject of ITU-R Rec. BS1286 [?]. A guide to existing recommendations that deal with audio quality assessment is provided in ITU-R Rec. BS1283 [?].

For video quality assessment, there exist many recommendations for assessing television picture quality. These recommendations are valid also for video in general. We can cite for example, ITU-T P.910 [?], etc. For multimedia quality assessment, the most used recommendation is ITU-T P.911 [?]. These recommendations provides some methods for non-interactive subjective assessment methods for evaluating the one-way overall audiovisual quality for multimedia applications. However, for interactive applications, there exist ITU-T P.920 [?], which provides the methods recommended for subjective interactive quality tests.

A. Source Signal for audiovisual Tests

ITU recommendations suggest that the duration of the source sequences should be about 10 seconds for audiovisual, and the length of 5 small different sentences for speech. In addition, it should be the actual duration of the sequence,

that means that the sequences cannot be obtained by repeating a shorter sequence. The termination of the sequences, should not cause an incomplete sentence or musical phrase. An initial and a final silent period or gray scene, not longer than 500 ms, can be used to make the sequence be more natural.

For the case of pair presentations, the references should have the best possible quality without any impairments. For audiovisual applications, speech and video should be perfectly synchronized.

B. Number of Subjects and their Selection

As stated in ITU recommendations, the number of subjects required to carry out the subjective quality test can vary from 4 to 40. Four is the absolute minimum for statistical reasons. The number of assessors needed depends upon the sensitivity and reliability of the test procedure adopted. The average number of subjects is about 15.

They should not be directly involved either in picture or audio quality evaluation as part of their work and should not be experienced assessors. Prior to a session, the observers should usually be screened for normal visual acuity or corrected-to-normal acuity and for normal color vision.

C. Instructions for Assessors

Before carrying out the experiments, some instructions should be given to the assessors. These instructions include the method of assessment, the types of impairment or quality factors likely to occur, the grading scale, the sequence and timing. These information should be explained and given to the subjects in a written form. The range and type of impairments should be presented in preliminary trials. Training trials may be given to subjects to familiarize them.

D. The Test Sessions

Following the ITU recommendations, overall subjective tests should be divided into multiple sessions and each session should not last more than 30 minutes. For each session, we should add several dummy sequences (about four or five) at the beginning. These sequences should not be taken into account in the calculation. Their aim is to be used as training samples for the subjects to learn how to give meaningful rates. Furthermore, the reliability of subjects can be qualitatively evaluated by checking their behavior when reference/reference pairs are given. In these cases, reliable subjects are expected to give evaluations very close to the maximum point in the quality scale. Hence, in these dummy samples, at least one of samples should be the reference. More advanced technique to check for the reliability of subjects is given in the next Section.

E. Subjective Quality Test Methods

In the ITU recommendations, there are many subjective quality test methods. Here we detail the methods that can be used with our method. Namely, the Absolute Category Rating (ACR), Degradation Category Rating (DCR) and The Double-Stimulus Continuous Quality-Scale (DSCQS) method are discussed here.

E.1 Absolute Category Rating (ACR)

The Absolute Category Rating method is a category judgment where the test sequences are presented one at a time and are rated independently on a category scale. (This method is also called Single Stimulus² Method.) In this method the subjects are asked to rate the quality of the presentation based on the level of the quality they have in their opinion for it after viewing or listening it. This phase is named the voting time. The time pattern for the stimulus presentation can be illustrated by Fig. 2. The voting time should be less than or equal to 10 seconds.

The five-level scale for rating overall quality is the most used scale, see Table I. If higher discriminative power is required, a nine-level scale may be used as shown in Table II. There is another variant of this scale which is the 11-point scale, depicted in Fig. 1. Finally, there is a general scale, which is the continuous quality scale that is depicted in Fig. 4 (which is rarely used with ACR method).

TABLE I
ITU 5-POINT QUALITY SCALE

Grading value	Estimated Quality
5	Excellent
4	Good
3	Fair
2	Poor
1	Bad

E.2 Degradation Category Rating (DCR)

The second suitable method for our tool, is the Degradation Category Rating (DCR). In this method the test sequences are presented in pairs. The first stimulus presented in each pair is always the source reference without any impairments. The second one is the same source but impaired by the test condition. This method is also called the Double Stimulus Impairment Scale (DSIS) method.

The time pattern for the stimulus presentation is depicted in Figure 3. The voting time should be less than or equal to 10 sec. In this case the subjects are asked to rate the impairment of the second stimulus in relation to the reference. The five-level scale for rating the impairment given in Table III is the most widely used one. However, all the quality scales used of ACR method can be used for DCR method but by replacing the quality adjectives by the corresponding

²The **stimulus** is the sequence presented to the subjects in a test.

TABLE II

ITU 9-POINT QUALITY SCALE

Grading value	Estimated Quality
9	Excellent
8	
7	Good
6	
5	Fair
4	
3	Poor
2	
1	Bad

TABLE III

ITU 5-POINT IMPAIRMENT SCALE

Grading value	Estimated impairment level
5	Imperceptible
4	Perceptible, but not annoying
3	Slightly annoying
2	Annoying
1	Very annoying

impairment adjectives. The method uses the impairment scale. As given in the ITU recommendations, the stability of the results is greater by using this scale.

F. The double-stimulus continuous quality-scale method

The third famous method is the double-stimulus continuous quality-scale (DSCQS) method. It is thought to be useful when it is not possible to provide test stimulus test conditions that exhibit the full range of quality. This method may be of special interest for our method to be used in some kind of networks guaranteeing certain level of QoS. For example,



Fig. 1. Eleven-point quality scale. 10 score for the sequence(s) that is identical to the reference one. Similarly, 0 score is for the sequence that has no similarity with the reference.

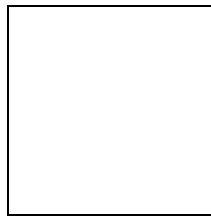


Fig. 2. Stimulus presentation timing in ACR method.



Fig. 3. Stimulus presentation timing in DCR method.

MOS greater than 3 on the 5-point quality scale.

The same as DCR method, the sequences are presented in pair: the reference and the impaired. Subjects are asked

to assess the quality of both (and not the impaired one with respect to the reference as in DCR). The unimpaired one is included to serve as a reference, but the observers are not told which is the reference sequence. In the series of tests, the position of the reference is changed randomly.

The subjects are asked to assess the overall sequence quality of each presentation by inserting a mark on a vertical scale (see Figure 5). The vertical scales are printed in pairs to accommodate the double presentation of each test sequence. The scales are continuous to avoid quantizing errors, but they are divided into five equal lengths which correspond to the normal ITU five-point quality scale. The associated terms categorizing the different levels are the same as those normally used. Here they are included for general guidance. Figure 5 shows a section of a typical score sheet.

G. Comparison of the methods

An important issue in choosing a test method is the fundamental difference between methods that use explicit references (e.g. DCR or DSCQS) and methods that do not use any explicit reference (e.g. ACR). The latter does not test fidelity. The former, on the other hand, should be used when testing the fidelity of transmission with respect to the source signal.

Thus, when it is important to check the fidelity with respect to the source signal, the DCR method should be used. Discrimination of imperceptible/perceptible impairment in the DCR scale supports this, as well as comparison with the reference quality. DSCQS in addition is used, as stated, in the cases when the quality range is not completely covered.

On the other hand, ACR is easy and fast to implement and the presentation of the stimuli is similar to that of the common use of the systems. Thus, ACR is well suited for qualification tests.

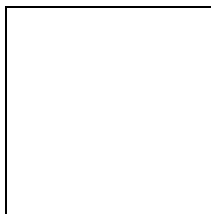


Fig. 4. A portion of quality rating form using continuous scales.



Fig. 5. A portion of quality rating form using continuous scales for DSCQS method.

IV. MOS CALCULATION

After performing any of the subjective methods, a range of integer values is given for each presentation. There will be variations in these distributions due to the differences in judgment between subjects. There are some analysis and calculations that should be performed on these data in order to remove the ratings of those subjects who could not conduct consistent scores. The most widely used reference that deal with this topic is the Recommendation of ITU-R BT.500-10, [?]. This is done to calculate the final mean opinion score (MOS) that we will use in the NN training and testing phases. In our method, we mean by a test condition, the resulting sample or sequence after distorting the original one (reference) by a selected set of values for the quality-affecting parameters.

For each of the test conditions, the mean and the 95% confidence interval of the statistical distribution of the assessment grades should be calculated. In addition, if the assessment was of change in impairment with a changing parameter value, as in our case, curve fitting techniques may be used.

The first step in the calculation of the mean score, \bar{u}_j , for each of the test conditions is given by:

$$\bar{u}_j = \frac{1}{N} \sum_{i=1}^N u_{ij}$$

where u_{ij} is the score of subject i , for test condition j , and N is the number of subjects. For each test condition, as well, a *confidence interval* should be calculated, which is derived from the standard deviation and the number of subjects. It is recommended to use the 95% confidence interval, given by:

$$[\bar{u}_j - \delta_j, \bar{u}_j + \delta_j]$$

where

$$\delta_j = 1.96 \frac{S_j}{\sqrt{N}}$$

and the estimated standard deviation, S_j , is given by

$$S_j = \sqrt{\sum_{i=1}^N \frac{(\bar{u}_j - u_{ij})^2}{N-1}}.$$

With a probability of 95%, the absolute value of the difference between the experimental mean score and the ‘true mean score (for a very high number of subjects) is then smaller than the 95% confidence interval.

A. Screening of Subjects

It must be ascertained whether this distribution of scores for test presentation is normal or not using the β_2 test (by calculating the kurtosis coefficient of the function, i.e. the ratio of the fourth order moment to the square of the second order moment). If β_2 is between 2 and 4, the distribution may be taken to be normal. For each test condition, the scores u_{ij} of each subject must be compared with the associated mean value, \bar{u}_j plus the associated standard deviation, S_j ,

times two (if normal) or times 20 (if non-normal), naming this summation P_j . Subject's scores are compared also to the associated mean value minus the same standard deviation times two or times 20, naming this quantity Q_j . Every time a subject's score is found above P_j a counter associated with each subject, P_i , is incremented. Similarly, every time a subject's score is found below Q_j a counter associated with each subject, Q_i , is incremented. Finally, the following two ratios must be calculated: $P_i + Q_i$ divided by the total number of scores from each subject for the whole session, and $P_i - Q_i$ divided by $P_i + Q_i$ as an absolute value. If the first ratio is greater than 5% and the second ratio is less than 30%, then subject i must be eliminated.

The process can be expressed mathematically as follows. For each test condition calculate the mean, the standard deviation, S_j , and kurtosis coefficient, β_{2j} , where β_{2j} is given by:

$$\beta_{2j} = \frac{m_{4j}}{(m_{2j})^2}, \quad \text{with} \quad m_{xj} = \frac{\sum_{i=1}^N (u_{ij} - \bar{u}_j)^x}{N}.$$

For each subject i , find P_i and Q_i (initialized to 0) using [?]

if $2 \leq \beta_{2j} \leq 4$, then:

if $u_{ij} \geq \bar{u}_j + 2S_j$ then $P_i = P_i + 1$;

if $u_{ij} \leq \bar{u}_j - 2S_j$ then $Q_i = Q_i + 1$;

else

if $u_{ij} \geq \bar{j} + \sqrt{20}S_j$ then $P_i = P_i + 1$;

if $u_{ij} \leq \bar{j} - \sqrt{20}S_j$ then $Q_i = Q_i + 1$.

If $(P_i + Q_i)/J > 0.05$ and $|(P_i - Q_i)/(P_i + Q_i)| < 0.3$ then reject subject i , where J is the number of test conditions.

For more details about this topic and the other methods of subjective tests see [?].

After eliminating the scores of those subjects who could not conduct coherent ratings using the above technique, the mean score should be recomputed using Eq. IV. This will constitute the MOS that we will use in our method.

It should be noted that there are two types of MOS. The first type referred as MOS (Mean Opinion Score), which uses the ACR (Absolute Category Rating) of the subjective quality tests. The other type is named Degradation MOS (DMOS) which uses the results from Degradation Category Rating (DCR).

V. AUTOMATIC MEASURING OF VIDEO QUALITY IN REAL TIME

In this Section, we summarize our proposal to evaluate video quality in an automatic way, in real time if necessary, and with results close to those that can be obtained from subjective tests. Then we briefly describe the fundamentals of RNN. After that, we compare ANN and RNN for the sake of our study. Finally, possible applications for our method are proposed.

A. Summary of the Neural Network approach

In this Section, we describe the overall steps to be followed in order to build a tool to automatically assess (and in real time if necessary) the quality of video streams transmitted over packet networks. The automatic evaluation is performed by a suitable trained Neural Network (NN).

First, we define a set of static information that will affect the general quality perception. We must choose the most effective quality-affecting parameters corresponding to the type of video application and to the network that will support the transmission. Then, for each parameter we must select the most frequent occurrences of its values and identify their ranges. For example, if the percentage loss rate is expected to vary from 0 to 10 %, then we may use 0, 1, 2, 5, and 10 % as typical values for this parameter. If we call *configuration* of the set of quality-affecting parameters, a set of values for each one, the total number of possible configurations is usually large. We must then select a part of this large cardinality set, which will be used as (part of) the input data of the NN in the learning phase.

To generate a video database composed of sequences corresponding to different configurations of the selected parameters (called “Distorted Database”), a simulation environment or a testbed must be implemented. This is used to send video sequences from the source to the destination and to control the underlying packet network. Every configuration in the defined input data must be mapped into the system composed of the network, the source and the receiver. For example, working with IP networks, the source controls the bit rate, the frame rate and the encoding algorithm, and it sends RTP video packets; the routers’ behavior contribute to the loss rate or the loss distribution, together with the traffic conditions in the network. The destination stores the transmitted video sequence and collects the corresponding values of the parameters. Then, by running the testbed or the using simulations, we produce and store a set of distorted video sequences along with the corresponding values of the parameters.

After completing the Distorted Database, a subjective quality test must be carried out. There are several subjective quality methods in the recommendations of the ITU-R [2]. We selected to use Degradation Category Rating (DCR), discussed in Section VI. A group of human subjects is then invited to evaluate the quality of the video sequences (i.e. every subject gives each video sequence a score from a predefined quality scale). The subjects should not establish any relation between the sequences and the corresponding parameters’ values.

The next step is to calculate the MOS values for all the video sequences. Based on the scores given by the human subjects, screening and statistical analysis may be carried out to remove the grading of the individuals suspected to give unreliable results [2]. After that, we store the MOS values and the corresponding parameters’ values in a second database (which we call the “Quality Database”).

In the third step, a suitable NN architecture and a training algorithm should be selected. The Quality Database is divided into two parts: one to train the NN and the other one to test its accuracy. The trained NN will then be able to evaluate the quality measure for any given values of the parameters.

To put this more formally, we build a set $\mathcal{S} = \{\sigma_1, \sigma_2, \dots, \sigma_S\}$ of video sequences that have encountered varied conditions when transmitted and that constitute the “training part” of the Quality Database. We also define a set $\mathcal{P} = \{\pi_1, \pi_2, \dots, \pi_P\}$ of parameters such as the bit rate of the source, the packet loss rate in the network, etc. Then, we denote by v_{ps} the value of parameter π_p in sequence σ_s , and by V the matrix $V = (v_{ps})$. For $s = 1, 2, \dots, S$, sequence σ_s receives the MOS evaluation $\mu_s \in [0, M]$ from the subjective test phase. The goal of the NN is to find a real function f having P real variables and with values in $[0, M]$, such that

- (i) for any sequence s , $f(v_{1s}, \dots, v_{Ps}) \approx \mu_s$,
- (ii) and such that for *any other* vector of parameter values (v_1, \dots, v_P) , $f(v_1, \dots, v_P)$ is close to the MOS that would receive any video sequence for which the selected parameters would have those specific values v_1, \dots, v_P .

The final tool is then composed of two modules: the first one collects the values of the selected quality-affecting parameters. The second one is the trained NN that will take the given values of the quality-affecting parameters and correspondingly computes the MOS quality score.

B. Random Neural Networks

The method we propose uses a new family of neural networks, the so called Random Neural Networks, recently invented by Erol Gelenbe in [16], [17], [18]. This choice was suggested by the success of this approach in many different areas [29], [30], [31], [32], ...

Gelenbe’s idea can be described as a merge between the classical Artificial Neural Networks (ANN) model and queuing networks. Since this tool is a novel one, let us describe here its main characteristics. RNN are, as ANN, composed of a set of interconnected neurons. These neurons exchange signals that travel instantaneously from neuron to neuron, and send and receive signals to and from the environment. Each neuron has a *potential* associated with, which is an integer (random) variable. The potential of neuron i at time t is denoted by $q_i(t)$. If the potential of neuron i is strictly positive, the neuron is *excited*; in that state, it randomly sends signals (to other neurons or to the environment), according to a Poisson process with rate r_i . Signals can be positive or negative. The probability that a signal sent by neuron i goes to neuron j as a positive one, is denoted by $p_{i,j}^+$, and as a negative one, by $p_{i,j}^-$; the signal goes to the environment (that is, it leaves the network) with probability d_i . So, if N is the number of neurons, we must have for all $i = 1, \dots, N$,

$$d_i + \sum_{j=1}^N (p_{i,j}^+ + p_{i,j}^-) = 1.$$

When a neuron receives a positive signal, either from another neuron or from the environment, its potential is increased by 1; if it receives a negative one, its potential decreases by 1 if it was strictly positive and it does not change if its value was 0. In the same way, when a neuron sends a signal, positive or negative, its potential is decreased by one unit (it was necessarily strictly positive since only excited neurons send signals)³. The flow of positive (resp. negative) signals

³In the general mathematical model, loops are allowed, that is, it is possible to have $p_{ii}^+ > 0$ or $p_{ii}^- > 0$. In our application, we set $p_{ii}^+ = p_{ii}^- = 0$.

arriving from the environment to neuron i (if any) is a Poisson process which rate is denoted by λ_i^+ (resp. λ_i^-). Finally, we make the usual independence assumptions between these arrival process, the processes composed of the signals sent by each neuron, etc.

The discovery of Gelenbe is that this model has a *product form* stationary solution. This is similar to the classical Jackson's result on open networks of queues. If process $\vec{q}(t) = (q_1(t), \dots, q_N(t))$ is ergodic (we will say that the network is *stable*), Gelenbe proved that

$$\lim_{t \rightarrow \infty} \Pr(\vec{q}(t) = (n_1, \dots, n_N)) = \prod_{i=1}^N (1 - \varrho_i) \varrho_i^{n_i} \quad (1)$$

where the ϱ_i s satisfy the following non-linear system of equations:

$$\text{for each node } i, \quad \varrho_i = \frac{T_i^+}{r_i + T_i^-}, \quad (2)$$

$$\text{for each node } i, \quad T_i^+ = \lambda_i^+ + \sum_{j=1}^N \varrho_j r_j p_{j,i}^+, \quad (3)$$

and

$$\text{for each node } i, \quad T_i^- = \lambda_i^- + \sum_{j=1}^N \varrho_j r_j p_{j,i}^-. \quad (4)$$

Relation (1) tells us that ϱ_i is the probability that, in equilibrium, neuron i is excited, that is,

$$\varrho_i = \lim_{t \rightarrow \infty} \Pr(q_i(t) > 0).$$

Observe that the non-linear system composed of equations (2), (3) and (4) has $3N$ equations and $3N$ unknowns (the ϱ_i s, the T_i^+ s and the T_i^- s). Relations (3) and (4) tell us that T_i^+ is the mean throughput of positive signals arriving to neuron i and that T_i^- is the corresponding mean throughput of negative signals (always in equilibrium). Finally, Gelenbe proved, first, that this non-linear system has a unique solution, and, second, that the stability condition of the network is equivalent to the fact that, for all node i , $\varrho_i < 1$.

Let us describe now the use of this model in statistical learning. Following previous applications of RNN, we fix the λ_i^- s to 0 (so, there is no negative signal arriving from outside). As a learning tool, the RNN will be seen as a black-box having N inputs and N outputs. The inputs are the rates of the incoming flows of positive signals arriving from outside, i.e. the λ_i^+ s. The output values are the ϱ_i s. In fact, in applications, most of the time some neurons do not receive signals from outside, which simply corresponds to fixing some λ_i^+ s to 0; in the same way, users often use as output only a subset of ϱ_i s.

At this point, let us assume that the number of neurons has been chosen, and that the topology of the network is selected; this means that we have selected the pairs of neurons that will exchange signals, without fixing the values of the rates r_i and the branching probabilities $p_{i,j}^+$ and $p_{i,j}^-$.

Our learning data is then composed of a set of K input-output pairs, which we will denote here by $\{(\vec{x}^{(k)}, \vec{y}^{(k)}), k = 1, \dots, K\}$, where $\vec{x}^{(k)} = (x_1^{(k)}, \dots, x_N^{(k)})$ and $\vec{y}^{(k)} = (y_1^{(k)}, \dots, y_N^{(k)})$. The goal of the learning process is to obtain values

for the remaining parameters of the RNN (the rates r_i and the branching probabilities $p_{i,j}^+$ and $p_{i,j}^-$) such that if, in the resulting RNN, we set $\lambda_i^+ = x_i^{(k)}$ for all i (and $\lambda_i^- = 0$), then, for all i , the steady-state occupation probability ϱ_i is close to $y_i^{(k)}$. This must hold for any value of $k \in \{1, \dots, K\}$.

To obtain this result, first of all, instead of working with rates and branching probabilities, the following variables are used:

$$w_{i,j}^+ = r_i p_{i,j}^+ \quad \text{and} \quad w_{i,j}^- = r_i p_{i,j}^-.$$

This means that $w_{i,j}^+$ (resp. $w_{i,j}^-$) is the mean throughput in equilibrium of positive (resp. negative) signals going from neuron i to neuron j . They are called *weights* by analogy to standard ANN. The learning algorithm proceeds then formally as follows. The set of weights in the network's topology is initialized to some arbitrary positive value, and then K iterations are performed which modify them. Let us call $w_{i,j}^{+(0)}$ and $w_{i,j}^{-(0)}$ the initial weights for the connection between i and j . Then, for $k = 1, \dots, K$, the set of weights at step k is computed from the set of weights at step $k-1$ using a *learning scheme* as usual with neural networks. More specifically, denote by $\mathcal{R}^{(k-1)}$ the network obtained after step $k-1$, defined by weights $w_{i,j}^{+(k-1)}$ and $w_{i,j}^{-(k-1)}$. When we set the input rates (external positive signals) in $\mathcal{R}^{(k-1)}$ to the $x_i^{(k)}$ s values, we obtain the steady-state occupations $\varrho_i^{(k)}$ s (assuming stability). The weights at step k are then defined by

$$w_{i,j}^{+(k)} = w_{i,j}^{+(k-1)} - \eta \sum_{l=1}^N c_l (\varrho_l^{(k)} - y_l^{(k)}) \frac{\partial \varrho_l}{\partial w_{i,j}^+}, \quad (5)$$

where the partial derivatives are evaluated for all h, m, n such that $\varrho_h = \varrho_h^{(k)}$ and $w_{m,n}^* = w_{m,n}^{*(k-1)}$. Weight $w_{i,j}^{-(k)}$ is defined in the same way. This is a gradient descent algorithm, corresponding to the minimization of the cost function

$$\frac{1}{2} \sum_{l=1}^N c_l (\varrho_l^{(k)} - y_l^{(k)})^2.$$

If some neuron l must not be considered in the output, we simply set $c_l = 0$. Once again, the relations between the output and the input parameters in the product form result allows to explicitly derive a calculation scheme for the partial derivatives. Instead of solving a non-linear system as (2), (3) and (4), it is shown in [17] that here we just have a linear system to be solved. When relation (5) is applied, it may happen that some value $w_{i,j}^{+(k)}$ or $w_{i,j}^{-(k)}$ is negative. Since this is not allowed in the model, the weight is set to 0 and it is no more concerned by the updating process (another possibility is to modify the η coefficient and apply the relation again; previous studies have been done using the first discussed solution, which we also adopt).

Once the K learning values have been used, the whole process is repeated several times, until some convergence conditions are satisfied. Remember that, ideally, we want to obtain a network able to give output $\vec{y}^{(k)}$ when the input is $\vec{x}^{(k)}$, for $k = 1, \dots, K$. The link with the notation in previous section is simply the following: K is the size of the training part of the Quality Database; for $k = 1, 2, \dots, K$, $x_i^{(k)} = v_i$ and for the (scalar) output, $y^{(k)} = \mu_k$.

As in most applications of ANN for learning purposes, we use a 3-level network structure: the set of neurons $\{1, \dots, N\}$

is partitioned into 3 subsets: the set of *input* nodes, the set of intermediate or *hidden* nodes and the set of *output* nodes. The input nodes receive (positive) signals from outside and don't send signals outside (that is, for each input node i , $\lambda_i^+ > 0$ and $d_i = 0$). For output nodes, the situation is the opposite: $\lambda_i^+ = 0$ and $d_i > 0$. The intermediate nodes are not directly connected to the environment; that is, for any hidden neuron i , we have $\lambda_i^+ = \lambda_i^- = d_i = 0$. Moreover, between the nodes on each level there are no transitions. Last, input neurons are only connected to hidden ones, and hidden neurons are only connected to output ones.

This is a typical structure for neural networks used as a learning tool. Moreover, RNN mathematical analysis (that is, solving the non-linear and linear systems) is considerably simplified in this case. In particular, it can easily be shown that the network is always stable. See again [17] or [18] for the details.

C. Comparison between ANN and RNN

In this subsection, we compare the two considered types of neural networks: Artificial Neural Networks (ANN) and Random Neural Networks (RNN), in the context of our specific problem. It is based on our experience with both models.

We used the Neural Networks MATLAB Toolbox when working with ANN, and a MATLAB package [33] for RNN. We observed that the ANN training process was relatively faster than that of RNN. However, during the run-time phase, RNN outperformed ANN in the total calculation time. This is because in the RNN's three-level architecture, the computation of the output ϱ_o for the unique output neuron o is done extremely fast: the non-linear system of equations (2), (3) and (4) allows, in this topology, to compute the ϱ_i s of the input layer directly, and of hidden layer neurons from the values for input layer ones. To be more specific, for each input neuron i we have

$$\varrho_i = \frac{\lambda_i^+}{r_i + \lambda_i^-},$$

(where, actually, we choose to set $\lambda_i^- = 0$), and for each hidden layer neuron h ,

$$\varrho_h = \frac{\sum_{\text{input neuron } i} \varrho_i w_{i,h}^+}{r_h + \sum_{\text{input neuron } i} \varrho_i w_{i,h}^-}.$$

The output of the black-box is then ϱ_o , given by

$$\varrho_o = \frac{\sum_{\text{hidden neuron } h} \varrho_h w_{h,o}^+}{r_o + \sum_{\text{hidden neuron } h} \varrho_h w_{h,o}^-}.$$

(The cost of computing the output ϱ_o is exactly $2IH + 3H + 1$ products (or divisions) and $H + 1$ sums, where I is the number of input neurons and H is the number of hidden ones.)

ANN's computations are slower because they involve nonlinear function calculations for each neuron in the architecture. This makes RNN particularly attractive for using them in contexts with real-time constraints, or for lightweight applications. This can be important in some kind of network applications, for example, in [34], an ANN packet-loss

predictor is proposed for real-time multimedia streams. The prediction precision is good, but the calculation time is much more than the next packet arrival time which makes the system useless unless for very powerful computers.

The most important feature of RNN we found for our problem is that it captures very well the mapping from parameters' values to the quality evaluation. This concerns also their ability to extrapolate in a coherent way for parameters' values out of the ranges used during the training phase. For instance, this led in [30] to build a zero-error channel decoder.

It is well known that the most common problems of ANN's learning are the over training and the sensitivity to the number of hidden neurons (the choice of the optimal number is difficult; one usually use heuristic methods for this). The over training problem makes the NN memorize the training patterns, but gives poor generalizations for new inputs. Moreover, if we can not identify some near-optimal number of hidden neurons, the performance may be bad for both the training set and the new inputs. Fig. 6 shows an example of an over-trained ANN network, where we can see irregularities, bad generalizations and bad capturing of the function mapping (see Subsection VIII-C for comparison with RNN).

We trained different architectures (varying the number of hidden neurons) for both ANN and RNN, with the same data (described in Section VII) and the same *mean square error*. Let us look, for instance, at the behavior of the quality as a function of the normalized bit rate BR (the 4 remaining variables were set to their most frequent observed values). In the database, BR varies between 0.15 to 0.7. In Fig. 7 and Fig. 8, we depict the ability of both networks to interpolate and extrapolate the results when BR varies from zero to its maximum value 1. Both figures show that RNN captures the mapping between the input and output variables, and that RNN is not very sensitive to the number of hidden neurons. While ANN gives quite different approximations for small changes in the size of the hidden layer.

Moreover, if the optimal ANN architecture could not be identified, its accuracy could be bad. Let us look now at the extreme values. If BR=0.0, the output should be around one, while for BR=1.0, the output should be between 8.5 and 9.0 on the y-axis. For the case of ANN, as shown in Fig. 8, when the number of hidden neurons changes from four (optimal experimentally) to five, the generalization is bad, specially when BR goes to zero, where the output is 3.2 instead of 1. This gives RNN a better ability to generalize.

D. Applications of Our Model

Let us briefly discuss here on the benefits from having a tool to automatically measure (and, if useful, in real time) video quality:

- In video conferencing and in the majority of video applications, at both end-user sides, such a tool can be used to monitor in real time the received video quality for control purposes. For example, changing the bit rate, using another codec, changing the frame rate, using some kind of FEC, changing the playback buffer size, etc. are possible decisions that can be taken to improve the quality or to maintain a certain level of the quality.

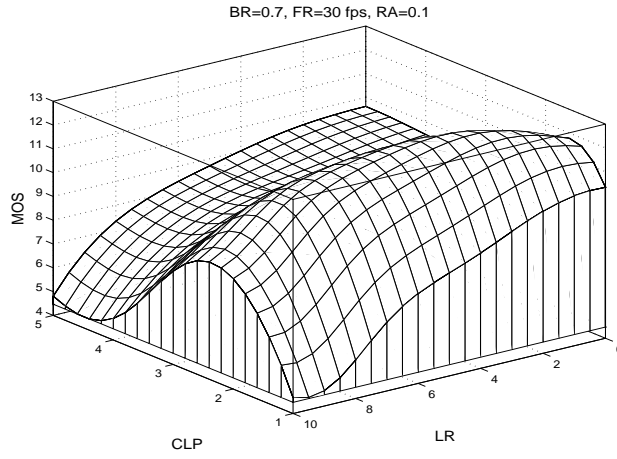


Fig. 6. Example of an over-trained ANN

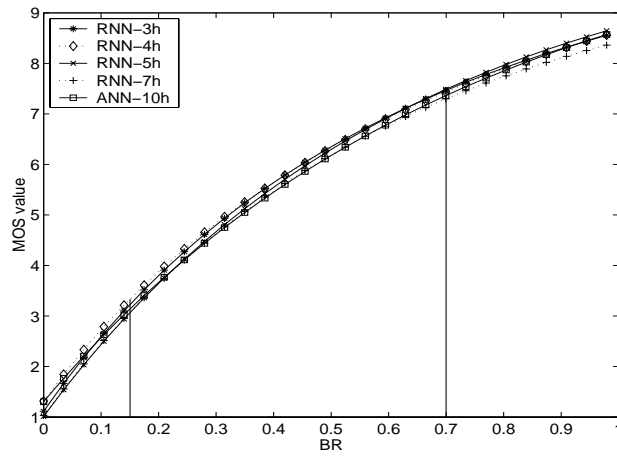


Fig. 7. Performance of RNN to interpolate and extrapolate for different number of hidden neurons

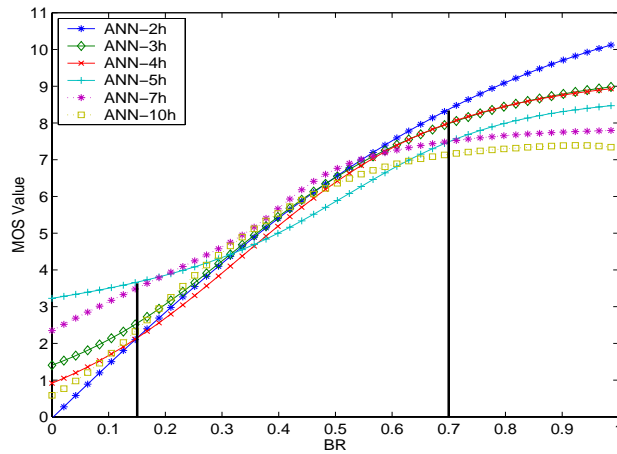


Fig. 8. Performance of ANN to interpolate and extrapolate for different number of hidden neurons

- Based on the ability of video quality measurement in real time, operators could use quality as a criteria for billing.
- The applications that transmit video on packet networks can use this tool to negotiate for the best configuration to

give the best quality.

- It can help also in the encoding process, as quality in encoders is also a way to “fit” the stream into the available global channel bandwidth. In video codecs using temporal compression (ex: MPEG, H.261, H.263...), a quality factor parameter is usually used to reduce the output stream bandwidth and to reduce, at the same time, the assessed quality (yet before any transmission). It will be even more interesting to have the history of all the important parameters (network and video) to compress the video signal rather than only the PSNR (Peak Signal to Noise Ratio) objective measure (which is poorer in quality and time consuming than our approach).

VI. SUBJECTIVE QUALITY AND MOS CALCULATIONS

As mentioned in the introduction, the main drawback of objective quality tests is that they do not correlate well with human perception. To evaluate the quality of video systems (codecs, telecommunications, television pictures, etc.), a subjective quality test is generally used. In this test, a group of human subjects is invited to judge the quality of the video sequence under different system conditions (distortions). There are several recommendations [2], [3] that specify strict conditions to be followed in order to carry out subjective tests. The main subjective quality methods are Degradation Category Rating (DCR), Pair Comparison (PC) and Absolute Category Rating (ACR). In our case, subjective tests are done using DCR.

In the DCR subjective quality test, a pair of video sequences is presented to each observer, one after the other. The observer must see the first one, which is not distorted by any impairment, and then the second one, which is the original signal distorted by some configuration of the set of chosen quality-affecting parameters. Fig. 9 shows the sequence and timing of presentations for this test. The time values come from the recommendation of the ITU-R [2].

The observer is asked to assess the overall quality of the distorted sequence with respect to the non-distorted one (the reference sequence), using a grade from one to nine corresponding to his/her mental measure of the quality associated with. It should be noted that there exist several quality scales. We chose this nine-grade one as a tradeoff between precision and dispersion of the subjective evaluations. Fig. 10 depicts the ITU-R nine-grade scale.

Following the ITU-R recommendations, overall subjective tests are divided into multiple sessions; each session should not last more than 30 minutes. Several dummy sequences are added in every session (about four or five) for training purposes. These sequences are not taken into account in the calculation. They are used to learn the observers how to give meaningful rates.

After performing any of the subjective methods, a range of integer values is given for each presentation. There will be variations in these distributions due to the differences in judgment between observers. For each of the test conditions, the mean score and a confidence interval are calculated. It is recommended to use a 95% confidence interval. A statistical analysis of the obtained results can also be carried out to remove the scores of observers suspected to give unreliable results (see Subsection VII-B). Details about this process and MOS calculations can be found in [2], [15].

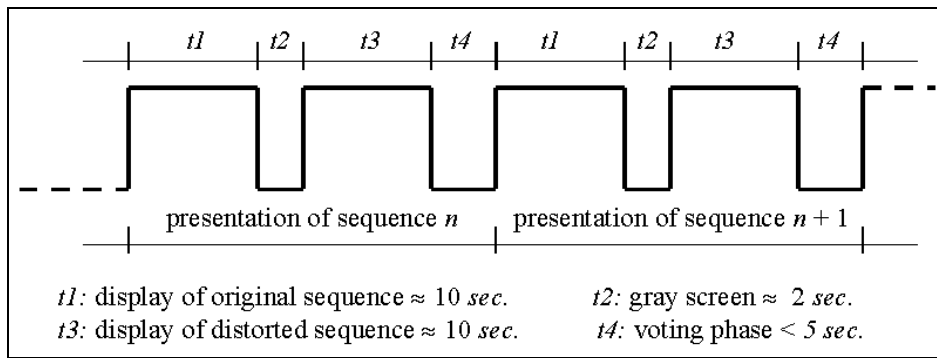


Fig. 9. Presentation structure of video sequences to the set of human observers in a DCR subjective quality test experiment.

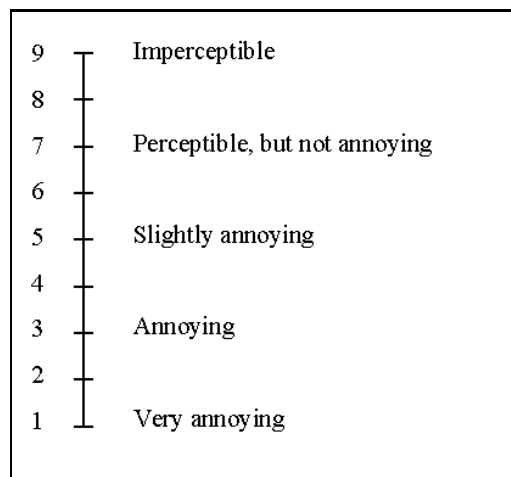


Fig. 10. The ITU-R nine-point quality scale

VII. PARAMETERS DESCRIPTION AND MOS TESTS

To generate the distorted video sequences, we used a tool that encodes a real-time video stream over IP networks into H263 format [35], simulates the packetization of the video stream, decodes the received stream, and allows us to simulate the network transmission conditions (packet loss process, etc.). The encoder can also be parameterized, in order to control the bit rate, the frame rate, the intra macro blocs refresh rate (i.e. it encodes the given macro bloc into intra mode rather than inter mode - this is done to make the stream more resistant to losses [36]), image format (QCIF, CIF, ...), etc. The packetization process is with conformance with RFC 2429 [37].

We used a standard video sequence called *stefan* to test the performance of H26x and MPEG4 codecs. It contains 300 frames encoded at 30 frames per sec., and lasts⁴ for 10 sec. The format of the encoded sequence is CIF (352 lines \times 288 pixels). The maximum allowed packet length is 536 bytes, in order to avoid the fragmentation of packets between routers. The use of one single sequence for the tests is to avoid any kind of semantic dependencies between the contents of the images and the evaluation made by the human subjects.

⁴This follows ITU recommendations, as usual in the area.

A. The Quality-Affecting Parameters

We present here the quality-affecting parameters that we consider having the highest impact on the quality:

- The Bit Rate (BR): this is the rate of the actual encoder's output. It is chosen to take four values (256, 512, 768 and 1024 Kbps.). It is known that not all the scenes compress the same ratio, for the same video quality, depending on the amount of redundancy in the scene (spatial and temporal), as well as image dimensions. All video encoders use a mixture of lossless and lossy compression techniques. Lossless compression does not degrade the quality, as the process is reversible. Lossy compression degrades the quality depending on the compression ratio needed by changing the quantization parameter of the Discrete Cosine Transform (DCT). For details about video encoding, see [35]. If the video is encoded only by lossless compression, the decoded video will have the same quality as the original one, provided that there is no other quality degradation. In our method, we normalize the encoder's output in the following way. If BR_{max} denotes the bit rate after the lossless compression process and BR_{out} is the final encoder's output bit rate, we select the scaled parameter $BR = BR_{out}/BR_{max}$. In our environments, we have $BR_{max} = 1430$ Kbps.
- The Frame Rate (FR): this is the number of frames per second. The original video sequence is encoded at 30 fps. The distorted sequence can be encoded at one of the following values: 6, 10, 15 and 30 fps. The encoder does this by dropping frames uniformly.
- The Loss Rate (LR): the simulator can drop packets randomly and uniformly to satisfy a given percentage loss rate. This parameter takes five values: 0, 1, 2, 4 and 10 %. This is because loss rates higher than 10 % drastically reduce video quality. In the networks where the LR is expected to be higher than this value, some kind of FEC [38] should be used to reduce the effect of losses.
- The number of Consecutively Lost Packets (CLP): we chose to drop packets in bursts of 1 to 5. The choice of 5 as upper limit for this parameter comes from real measurements that we performed before [27]. Regarding the loss model (LR and CLP), we used the same one as [7].
- The ratio of the encoded intra macro-blocs to inter macro-blocs (RA): the encoder can change the refresh rate of the intra macro-blocs in order to make the encoded sequence more or less sensitive to packet losses [36]. This parameter takes values that vary between 0.05 and 0.50 depending on the BR and the FR. We selected five values for it.

The delay and the delay variation are indirectly considered: they are included in the LR parameter. Indeed, it is known that if a de-jittering mechanism with a strict playback buffer length is used, then all the packets arriving after a predefined threshold are considered as lost. So, in this way, all delays and delay variations are mapped into loss. This is confirmed in [26].

The whole set of configurations of these 5 parameters has 2000 elements. To build the Distorted Database, we first selected a default value for each parameter (the most frequent observed value). Then, for each possible selection of 3 parameters among the 5 (10 possibilities), we set the 3 parameters to their defaults values, and we built all the possible

combinations of the remaining 2. In this way, a list of configurations was defined, containing many duplications. Once these duplications removed, we obtained a set of 94 different configurations.

B. Subjective Quality Test and the RNN Architecture

The subjective quality test is with conformance to the method Degradation Category Rating (DCR), with a quality scale consisting of 9 points (see Section VI). We divided the test into two sessions, and added 5 distorted sequences to the first session and 4 to the second one. These nine sequences were not considered in the MOS calculation since they are used as a training phase for the human subjects. At the same time, they are used to verify how much reliable is the person carrying out the test, as they are replicated from the real 94 samples.

We invited 20 persons to perform the subjective tests. After that, a statistical analysis of the results was performed; as a consequence, we discarded the notes of two subjects. Then, we calculated the MOS scores and the corresponding parameters' values used to train the RNN. We used as architecture of RNN, the three-layer feedforward consisting of five neurons in the input layer (which corresponds to the five chosen parameters), one output neuron (corresponding to the MOS scores), and 5 neurons in the hidden layer. Indeed, we tested several other sizes for this layer, and we found that 5 hidden neurons were sufficient for our purposes.

VIII. PARAMETERS IMPACT ON VIDEO QUALITY

In this section, we study the effect of the mentioned five parameters on video quality. As it is impossible to visualize the variation of the quality as a function of all the five parameters, we chose to visualize on a set of 3-D figures in which we varied two parameters and kept the other three fixed. The MOS value is computed as a function of the values of all the five parameters by using the RNN tool described in the previous Section. It should be noted that the axis of these figures are rotated in such a way to give the best visualization of each figure. This means that the orientation of some axis can vary from figure to figure. In addition, the scale of quality axis (or Z-axis) is variable from one figure to another to show some useful features of the given figure. In the following subsections, we analyze the effect of each parameter on the quality and we also discuss their main combined effects.

A. Bit Rate (BR)

As shown in Fig. 11 through Fig. 14, the parameter BR exhibits a great influence on the quality. From Fig. 11, MOS values, obtained by the trained RNN, vary from 3 to 6.5 for (FR=6 fps, LR=0 %, CLP=2, and RA=0.1) and from 3 to 7.8 when the FR goes up to 30 fps. This is for a variation of BR from 0.15 to 0.8. Its effect is comparable to the effect of LR as shown in Fig. 12. The improvement of the quality is significant for lower LR and decreases for higher LR. The quality goes from 2.2 to 8.0 for zero loss, but only about 1.1 of absolute improvement for 10% LR. As shown in Fig. 13 and Fig. 14, for higher values of BR, there is no effect of the variations of CLP and RA parameters on the quality, however; their effects increase and become benefic for lower values of BR.

When the encoder has to decrease the BR of the given stream, it increases the quantization parameter in order to further compress the original image. This process increases the artifacts of the output stream, and hence increases the distortion which becomes more noticeable by the human as the BR decreases. This parameter is studied in [8] and [1] however the interaction between other parameters is not shown.

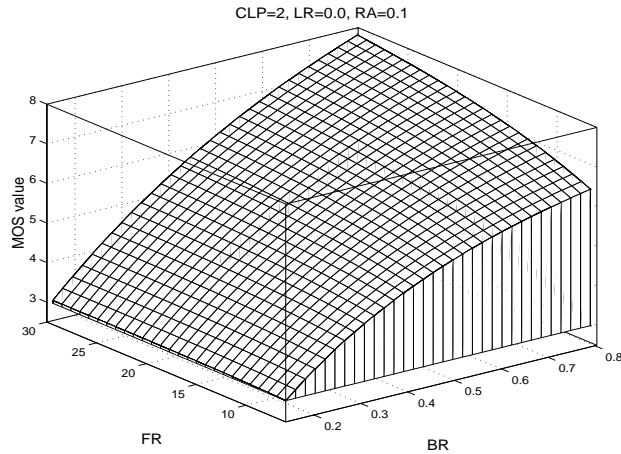


Fig. 11. The impact of BR and FR on video quality.

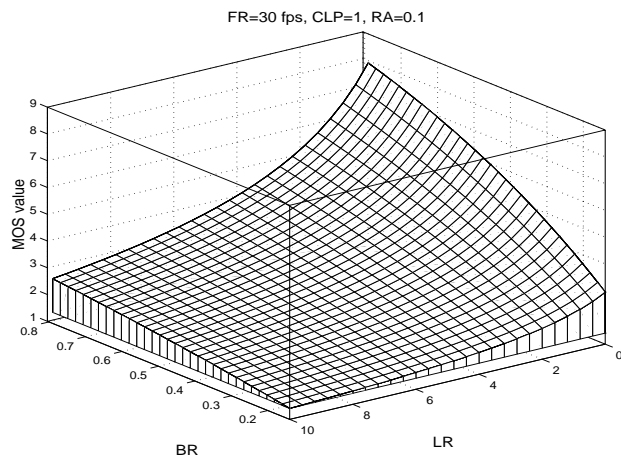


Fig. 12. The impact of BR and LR on video quality.

B. Frame Rate (FR)

The effect of this parameter on the quality is not as significant as in the BR or LR (see below) cases. This is clearly shown in Figs. 11, 15, 16 and 17. For (BR=0.8, LR=0 %, CLP=2, and RA=0.1), an improvement from 5.1 up to 7.3 is achieved for a wide range variation going from 6 to 30 fps, but no impact on the quality for BR=0.15, as depicted in Fig. 11. We can note also that the enhancement of the quality for FR greater than 16 is negligible, as already observed in previous works [11]. As the BR decreases, the effect of FR becomes smaller, see Fig. 21.

The improvement of the quality when the LR changes from zero to 10 % (MOS value varies from 6.2 up to 7.5 for

FR=30 fps) decreases until it becomes negligible as shown in Fig. 15. While from Fig. 17, we can see that we can get constant relative improvements of the quality whatever the value of RA (from 5.8 to 6.9 for RA=0.05 and from 6.0 to 7.1 for RA=0.5 - this is for BR=0.6, CLP=2 and LR=0 %). With respect to CLP parameter we have similar results (see Fig. 16).

These results may seem surprising, but they have been observed in different previous studies. For example, experimental results showed that for FR larger than 16 fps, the viewer is not so sensitive to changes in FR values [11]. The work done in [9] clearly showed that the enhancement of the quality for a wide range variation from 6 to 25 fps is really small. Both studies used subjective quality tests.

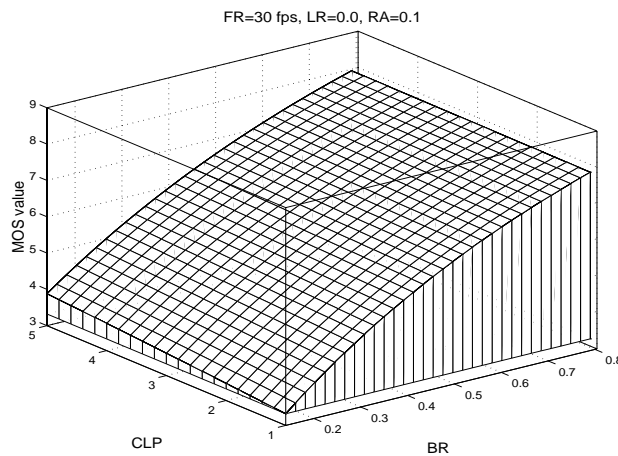


Fig. 13. The impact of BR and CLP on video quality.

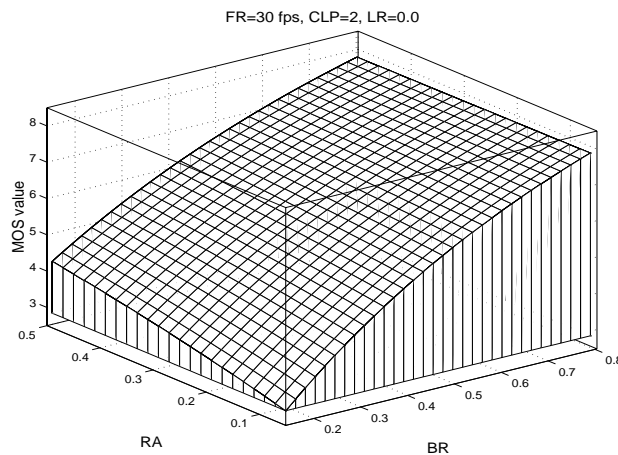


Fig. 14. The impact of BR and RA on video quality.

C. Loss Rate (LR)

The effect of LR on video quality was the main goal of several previous studies as it is an important parameter [21], [9], [7]. As shown in Fig. 12, the quality drastically decreases when the LR increases from 0 up to 10 %. But the absolute

decrease of the quality depends on the other parameters' values: about 1.0 for BR=0.15, while 6.9 for BR=0.8. As depicted in Fig. 15, the absolute deterioration is about 3.8 for 6 fps, but 5.0 for 30 fps. The variation of the quality is about 3.8 for RA=0.5 and 3.0 for RA=0.05, when BR=0.35, FR=30 fps and CLP=2, as shown in Fig. 19.

As we can see from Fig. 18, for zero loss, there is no change on the quality for any value of CLP. This may seem obvious, but it shows that RNN captures with great confidence and reliability the characteristics of the quality as a function of the parameters. It may be thought that this is due to large values of BR and FR, but we varied all these parameters and we got the same behavior as shown in Fig. 21. In the case of using ANN, it is not easy to get this result (see Fig. 6).

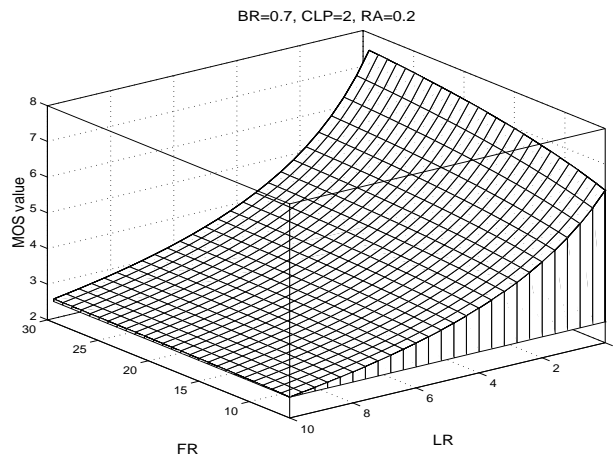


Fig. 15. The impact of FR and LR on video quality.

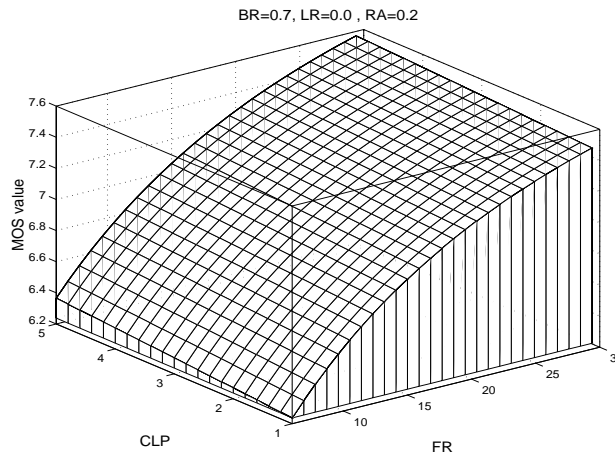


Fig. 16. The impact of FR and CLP on video quality.

D. The number of Consecutively Lost Packet (CLP)

As the result obtained for this parameter may seem strange, let us clarify its effects. When keeping all other parameters constant and increasing the value of CLP, the distance between any two consecutive loss occurrences increases. This has

two consequences: loss occurrences decrease and consequently the deterioration of frames becomes smaller (i.e. smaller number of frames partially distorted by loss). As it is shown from our previous analysis of FR, the eye is less sensitive to higher values of FR. Moreover, for each lost packet the past macro-blocs (a portion of the image) are still shown on the screen as a kind of error resilience. Hence, larger values of CLP may introduce deterioration in smaller frames. This is equivalent to smaller LR values but slightly lower FR values. As it was previously shown that the effect of LR is much more effective than that of FR, the effect of higher values of CLP is benefic to the quality. This result is with agreement with that obtained in [7], although the authors did not explain why this is happening.

Starting by Fig. 13, and for larger values of BR, the effect of CLP is negligible, but as the BR decreases the effect of CLP increases up to 0.75 for BR=0.15. As shown in Fig. 16, the effect of CLP is small, only an improvement of 0.2, but this is due to the fixed values of the other parameters (BR=0.7, LR=0 and RA=0.2). For smaller values of BR and larger values of LR, the effect is much more noticeable.

From Fig. 20, we can see that both CLP and RA has small effect on the quality (note the quality scale). In conclusion, CLP effect is important and comparable to the FR effect, but it is more benefic in the case of poor conditions (i.e. smaller BR or more LR), while the FR is benefic for good conditions (higher BR and lower LR).

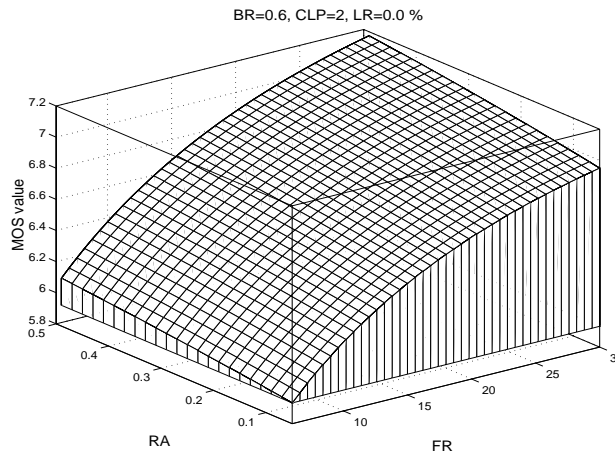


Fig. 17. The impact of FR and RA on video quality.

E. Intra-to-Inter Ratio (RA)

From the definition of RA, it is the ratio of the encoded intra macro-blocs to those encoded as inter, so it is expected that the effect of this parameter is benefic for increase of its value. This is clearly shown in Figs. 14, 17, 19 and 20. Its effect is more than the FR parameter, and more interestingly for lower values of the BR parameter. This is shown in Fig. 14, where for BR=0.15 the quality increased from 2.8 up to 4.2. This improvement is relatively important when the network bandwidth is small. But for larger BR, the effect is negligible.

The two parameters LR and RA are related, as shown in Fig. 19. For smaller LR, there is great improvement on the

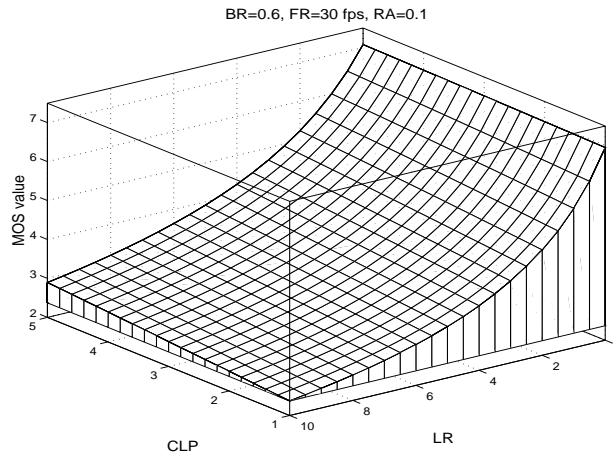


Fig. 18. The impact of LR and CLP on video quality.

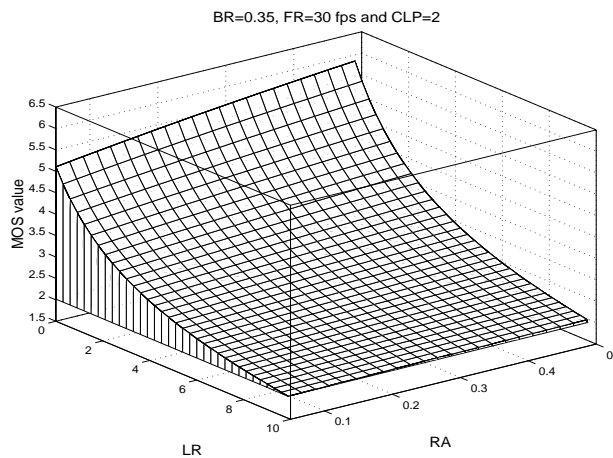


Fig. 19. The impact of LR and RA on video quality.

quality for the increase in RA. But this improvement vanishes when the LR increases. That means that we can get better video quality when the available channel bandwidth (BR) is small and for smaller values of LR by increasing the value of RA. In such a case, the effect of RA on the quality is more beneficial than that of FR, as shown in Fig. 21.

IX. CONCLUSIONS

In this paper, we first propose a methodology that is able to automatically quantify the quality of a video stream at the receiver side, after its transport over a packet network. The main properties of our technique are the following: (i) the evaluation can be done in real time, (ii) the value assigned to the quality of the flow is close to the value that could be obtained from a subjective test, because it is done by a neural network that was trained to behave as the average of a set of human observers of the video stream. Once the neural network is trained, the tool continuously receives the values of the selected parameters (from measurements) to perform its evaluation.

Second, we used our technique to study the impact on video quality of some of the most important parameters

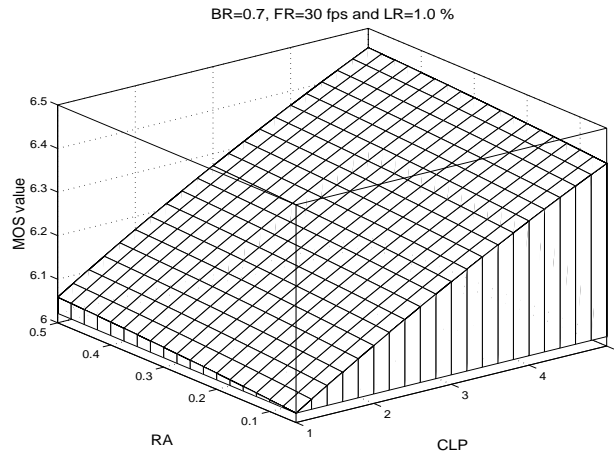


Fig. 20. The impact of CLP and RA on video quality.

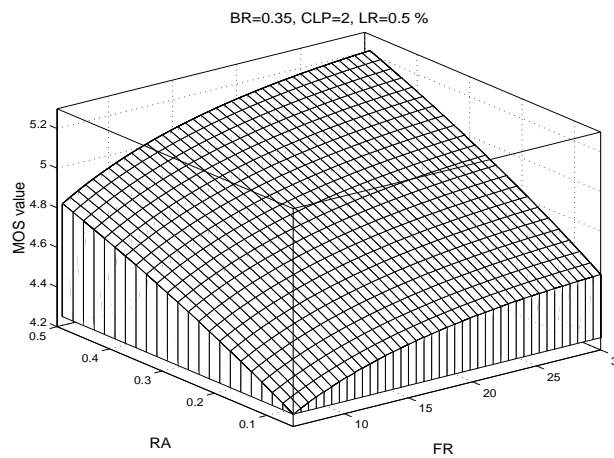


Fig. 21. RA is more benefic than FR for lower values of BR.

supposed to have an impact on it, namely, the stream bit rate, the scene frame rate, the network loss rate, the burst loss size and the ratio of the encoded intra to inter macro-blocs for H263 codecs. As far as we know, there is no previous study of the combined effect of several parameters on video quality. The goal of this analysis is to help in the understanding of the behavior of real-time video streams transmitted over best-effort networks. This may be used, for instance, in developing control mechanisms allowing the delivery of the best possible video quality given the current network state.

To achieve these goals, we improved our previous proposals (in both audio and video transmission) using a new and efficient neural model called Random Neural Network, which behaves significantly better than the standard one.

Although we based our study on IP networks, our approach can be followed for ATM and wireless technologies as well (the specific type of packet technology has no effect on the relevance of the method we discuss here). Some future research directions for this work include the study of other codecs like MPEG2/4, the analysis of the effect of audio and video synchronization. Concerning the set of selected parameters, we also intend to explore the usefulness of more sophisticated models on the different aspects that we analyze here.

REFERENCES

- [1] O. Olsson, M. Stoppiana, and J. Baina, "Objective methods for assessment of video quality: state of the art," *IEEE Transactions on Broadcasting*, vol. 43, no. 4, December 1997.
- [2] Rec. ITU-R BT.500-10, "Methodology for the subjective assessment of the quality of television pictures," in *International Telecommunication Union*, March 2000.
- [3] Rec. ITU-T P.910, "Subjective video quality assessment methods for multimedia applications," in *International Telecommunication Union*, September 1999.
- [4] S.D. Voran and S. Wolf, "The development and evaluation of an objective quality assessment system that emulates human viewing panels," in *International Broadcasting Convention*, 1992.
- [5] C. Jones and D.J. Atkinson, "Development of opinion-based audiovisual quality models for desktop video-teleconferencing," in *International Workshop on Quality of Service*, 1998.
- [6] K.T. Tan and M. Ghanbari, "A combinational automated MPEG video quality assessment model," in *Proc. of Image Processing and its Application Conference*, 1999.
- [7] D. Hands and M. Wilkins, "A study of the impact of network loss and burst size on video streaming quality and acceptability," in *Interactive Distributed Multimedia Systems and Telecommunication Services Workshop*, 1999.
- [8] H.R. Wu, T. Ferguson, and B. Qiu, "Digital video quality evaluation using quantitative quality metrics," in *Proc. of the 4th International Conference on Signal Processing*, 1998.
- [9] G. Ghinea and J.P. Thomas, "QoS impact on user perception and understanding of multimedia video clips," in *Proc. ACM Multimedia 98*, 1998.
- [10] M Claypool and J. Tanner, "The effects of jitter on the perceptual quality of video," in *ACM Multimedia 99*, 1999.
- [11] K. Shuaib, T. Saadawi, M. Lee, and B. Basch, "A de-jittering scheme for the transport of MPEG-4 and MPEG-2 video over ATM," *IEEE Military Communications Conference Proc., MILCOM*, 1999.
- [12] H.R. Wu, C. Lambrecht, M. Yuen, and B. Qiu, "Quantitative quality and impairment metrics for digitally coded images and image sequences," in *Proc. of Australian Telecommunication Networks and Applications Conference*, December 1996.
- [13] M. Caramma, R. Lancini, and M. Marconi, "Subjective quality evaluation of video sequences by using motion information," in *Proc. of International Conference on Image Processing, ICIP'99*, October 1999.
- [14] A. Watson and M.A. Sasse, "Measuring perceived quality of speech and video in multimedia conferencing applications," in *Proc. of ACM Multimedia'98*, September 1998, pp. 55–60.
- [15] S. Mohamed, G. Rubino, H. Afifi, and F. Cervantes, "Real-time video quality assessment in packet networks: A neural network model," in *Proc. International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA'01)*, June 2001.
- [16] E. Gelenbe, "Random neural networks with negative and positive signals and product form solution," in *Neural Computation*, 1989, vol. 1, pp. 502–511.
- [17] E. Gelenbe, "Learning in the recurrent random neural network," in *Neural Computation*, 1993, vol. 5, pp. 154–511.
- [18] E. Gelenbe, "Stability of the random neural network model," in *Neural Computation*, 1990, vol. 2, pp. 239–247.
- [19] A. Vahedian, M.R. Frater, and J.F. Arnold, "Impact of audio on subjective assessment of video quality," in *Proc. of International Conference on Image Processing*, 1999.
- [20] A. Watson and M.A. Sasse, "Evaluating audio and video quality in low-cost multimedia conferencing systems," *ACM Interacting with Computers Journal*, vol. 8, no. 3, pp. 255–275, 1996.
- [21] J.M. Boyce and R.D. Gaglianella, "Packet loss effects on MPEG video sent over the public internet," in *Proc. ACM Multimedia'98*, 1998.
- [22] S. Lavington, N. Dewhurst, , and M. Ghanbari, "The performance of layered video over an IP network," in *Packet Video Workshop*, May 2000.

- [23] D. Qiao and F.Y. Zheng, "Dynamic bit-rate estimation and control for constant-quality communication of video," in *Proc. of the 3rd World of Intelligent Control and Automation*, 2000.
- [24] Y. Inazumi, Y. Horita, K. Kotani, and T. Murai, "Quality evaluation method considering time transition of coded video quality," in *Proc. of International Conference on Image Processing*, 1999.
- [25] A. Pessoa, A. Falcao, R. Nishihara, A. Silva, and R. Lotufo, "Video quality assessment using objective parameters based on image segmentation," *SMPTE Journal*, December 1999.
- [26] R. Zhang, R. Regunathan, and K. Rose, "Video coding with optimal Inter/Intra-mode switching for packet-loss resilience," *IEEE Journal on Selected Area in Communications*, vol. 18, no. 6, June 2000.
- [27] S. Mohamed, F. Cervantes, and H. Afifi, "Audio quality assessment in packet networks: an inter-subjective neural network model," in *Proc. of the 15th IEEE International Conference on Information Networking (ICOIN-15)*, January 2001.
- [28] S. Mohamed, F. Cervantes, and H. Afifi, "Integrating networks measurements and speech quality subjective scores for control purposes," in *Proc. of the IEEE INFOCOM'01*, April 2001.
- [29] C. Cramer, E. Gelenbe, and P. Gelenbe, "Image and video compression," in *IEEE Potentials*, 1998.
- [30] A.H. Abdelbaki and E. Gelenbe, "Random neural network decoder for error correcting codes," in *International Joint Conference on Neural Networks*, 1999, vol. 5.
- [31] E. Gelenbe, H. Bakircioglu, and T. Kocak, "Image processing with the random neural network," in *Proc. of the IEEE Digital Signal Processing Conference*, June 1997.
- [32] H. Bakircioglu and T. Kocak, "Survey of random neural network applications," *European Journal of Operational Research* 126, 2000.
- [33] H. Abdelbaki, "Random neural network simulator," <ftp://ftp.mathworks.com/pub/contrib/v5/nnet/rnnsimv2/>.
- [34] S. Lavington, H. Hagrass, and N. Dewhurst, "Using a MLP to predict packet loss during real-time video transmission," Tech. Rep. CSM329, Univ. Of Essex, UK, July 1999.
- [35] ITU-T Recommendation H.263, "Video coding for low bit rate communication," in *International Telecommunication Union*, February 1998.
- [36] F. Le Leannec and C. Guillemot, "Packet loss resilient H.263+ compliant video coding," in *International Conference on Image Processing*, September 2000.
- [37] RFC 2429, "RTP payload format for the 1998 version of ITU-T Rec. H.263 video H.263+," in *IETF*, October 1998.
- [38] W. Tan and A. Zakhor, "Multicast transmission of scalable video using receiver-driven hierarchical FEC," in *Packet Video Workshop 99*, April 1999.

Samir Mohamed obtained a B.Sc. degree in computer and control systems, from the Faculty of Engineering at the university of Assuit, Egypt, in May 1994. He worked as teaching assistant in the Faculty of Engineering at Asswan from 1995 to 1997. He obtained his M.Sc. degree in computer science from the university of Rennes I, France, in 1998. He is now a PhD student at INRIA/IRISA, Rennes, France. His research interests are in audio and video quality assessment in packet networks and rate based control mechanisms.

Gerardo Rubino is a Senior Researcher at INRIA, in the Center of Rennes, France. He received his PhD in Computer Science from the university of Rennes 1, France, in 1985. Since then, he performs research work for INRIA. Since 1995, he also works at the Networks and Multimedia Services Department at the ENST Bretagne engineering school, as the responsible for the research activities in that group. His main interests are in the quantitative analysis of computer and communications systems, mainly using probabilistic models. He is the head of a team at INRIA working in the design and analysis of communications networks. He currently serves in the Editorial Board of the Naval Research Logistics