



**INTERNATIONAL ELECTROTECHNICAL COMMISSION**

**TECHNICAL COMMITTEE NO. 100: AUDIO, VIDEO AND MULTIMEDIA SYSTEMS AND EQUIPMENT**

**PROJECT TEAM NO. 62251: MULTIMEDIA SYSTEMS AND EQUIPMENT – QUALITY ASSESSMENT – AUDIO-VIDEO CONFERENCING SYSTEMS**

2<sup>nd</sup> Working Draft for Project 62251

-----

The attached working draft is prepared by a small group in Japan composed of some of PT 62251 members, domestic members and PL, including Vice Chairman of ITU-R WP6Q, after the virtual meeting immediately following the voting on NP initiated by the Japanese National Committee.

Please note that the title of the working document has been changed as proposed at the virtual meeting and reported to TC 100/AGM meeting in May 2001 in Brussels;

**“Multimedia systems and equipment – Quality assessment – Audio-video communication systems”**

This document is of attention of all PT 62251 members in advance to the first physical meeting of PT 62251 on 2001-10-18 in Firenze (Italy).

Any constructive comments with change proposals and contributions are welcome to improve this document. They should be posted to the mailing list or sent directly to Project Leader at the address below.

This document together with received contributions will be discussed under the agenda item 5 at the forthcoming PT 62251 meeting in Firenze.

-----  
Project Leader for IEC TC 100/PT 62251  
Hiroaki Ikeda  
lkeda@hike.tu.chiba-u.ac.jp

## CONTENTS

	Page
1 Scope.....	4
2 References.....	4
3 Terms and definitions .....	5
4 Configuration for quality assessment.....	5
4.1 Input and output channels .....	5
4.2 Points of input and output terminals .....	5
5 Video quality .....	6
5.1 End-to-end tone reproduction .....	6
5.1.1 Item to be assessed .....	6
5.1.2 Method of assessment.....	6
5.1.3 Form of reporting assessment result.....	7
5.2 End-to-end colour reproduction .....	8
5.2.1 Item to be assessed .....	8
5.2.2 Method of assessment.....	8
5.2.3 Form of reporting assessment result.....	9
5.3 Peak-signal to noise ratio (PSNR) .....	10
5.3.1 Item to be assessed .....	10
5.3.2 Method of assessment.....	10
5.3.3 Form of reporting assessment results .....	11
6 Audio quality .....	15
6.1 Perceived audio quality with full-reference signals.....	15
6.1.1 Item to be assessed .....	15
6.1.2 Justification .....	15
6.1.3 Method of assessment and algorithm of PEAQ .....	15
6.1.4 Form of reporting assessment results .....	16
6.2 Sampling rate and quantization resolution .....	17
6.2.1 Item to be assessed .....	17
6.2.2 Method of assessment.....	17
6.2.3 Form of reporting assessment results .....	17
6.3 Delay .....	18
6.3.1 Item to be assessed .....	18
6.3.2 Method of assessment.....	18
6.3.3 Form of reporting assessment result.....	18
7 Total quality .....	19
7.1 Synchronization of audio and video (lip sync) .....	19
7.1.1 Item to be assessed .....	19
7.1.2 Method of assessment.....	19
7.1.3 Form of reporting assessment results .....	20
7.2 Scalability .....	20
7.2.1 Item to be assessed .....	20
7.2.2 Method of assessment.....	20
7.2.3 Form of reporting assessment results .....	20
7.3 Overall quality .....	20
7.3.1 Item to be assessed .....	20

- 7.3.2 Method of assessment..... 20
- 7.3.3 Form of reporting assessment results ..... 20
- A.1 Introduction ..... 21
- A.2 Test sources and hypothetical deterioration..... 21
- Annex B (informative) PEAQ objective measurement method outline ..... 25
  - B.1 Basic concept of the PEAQ measurement algorithm ..... 25
  - B.2 Basic version..... 26
  - B.3 Advanced version..... 27
  - B.4 Output value of PEAQ method..... 28
  - B.5 Performance of PEAQ measurement method..... 28
- Bibliography ..... 29

## INTERNATIONAL ELECTROTECHNICAL COMMISSION

**MULTIMEDIA SYSTEMS AND EQUIPMENT – QUALITY ASSESSMENT –  
AUDIO-VIDEO COMMUNICATION SYSTEMS**

## FOREWORD

- 1) The IEC (International Electrotechnical Commission) is a worldwide organization for standardization comprising all national electrotechnical committees (IEC National Committees). The object of the IEC is to promote international co-operation on all questions concerning standardization in the electrical and electronic fields. To this end and in addition to other activities, the IEC publishes International Standards. Their preparation is entrusted to technical committees; any IEC National Committee interested in the subject dealt with may participate in this preparatory work. International, governmental and non-governmental organizations liaising with the IEC also participate in this preparation. The IEC collaborates closely with the International Organization for Standardization (ISO) in accordance with conditions determined by agreement between the two organizations.
- 2) The formal decisions or agreements of the IEC on technical matters express, as nearly as possible, an international consensus of opinion on the relevant subjects since each technical committee has representation from all interested National Committees.
- 3) The documents produced have the form of recommendations for international use and are published in the form of standards, technical specifications, technical reports or guides and they are accepted by the National Committees in that sense.
- 4) In order to promote international unification, IEC National Committees undertake to apply IEC International Standards transparently to the maximum extent possible in their national and regional standards. Any divergence between the IEC Standard and the corresponding national or regional standard shall be clearly indicated in the latter.
- 5) The IEC provides no marking procedure to indicate its approval and cannot be rendered responsible for any equipment declared to be in conformity with one of its standards.
- 6) Attention is drawn to the possibility that some of the elements of this technical report may be the subject of patent rights. The IEC shall not be held responsible for identifying any or all such patent rights.

The main task of IEC technical committees is to prepare International Standards. However, a technical committee may propose the publication of a technical report when it has collected data of a different kind from that which is normally published as an International Standard, for example "state of the art".

Technical reports do not necessarily have to be reviewed until the data they provide are considered to be no longer valid or useful by the maintenance team.

IEC 62251, which is a technical report, has been prepared by IEC technical committee 100: Audio, Video and Multimedia Systems and Equipment.

The text of this technical report is based on the following documents:

Enquiry draft	Report on voting
XX/XX/CDV	XX/XX/RVC

Full information on the voting for the approval of this technical report can be found in the report on voting indicated in the above table.

This publication has been drafted in accordance with the ISO/IEC Directives, Part 2.

This document which is purely informative is not to be regarded as an International Standard.

## MULTIMEDIA SYSTEMS AND EQUIPMENT – QUALITY ASSESSMENT – AUDIO-VIDEO COMMUNICATION SYSTEMS

### 1 Scope

This Technical Report specifies items to be measured by objective methods, methods of measurement together with measuring conditions, processing of the measured data and forms to report acquired information for assessment of end-to-end quality of audio-video communication systems over digital networks. The measurements are supposed to be conducted in a double-ended and a full reference. The systems are assumed to have electrical interface channels at the input and at the output of audio-video signals for objective assessment.

The extension for systems that do not have such channels is left for further study.

### 2 References

The following normative documents contain provisions which, through reference in this text, constitute provisions of this International Standard. For dated references, subsequent amendments to, or revisions of, any of these publications do not apply. However, parties to agreements based on this International Standard are encouraged to investigate the possibility of applying the most recent editions of the normative documents indicated below. For undated references, the latest edition of the normative document referred to applies. Members of IEC and ISO maintain registers of currently valid International Standards.

IEC 61146-1: 1994, Video cameras (PAL/SECAM/NTSC) – Methods of measurement – Part 1: Non-broadcast single-sensor cameras.

IEC 61146-2: 1997, Video cameras (PAL/SECAM/NTSC) – Methods of measurement – Part 2: Two- and three-sensor professional cameras.

IEC 61966-2-1: 1999, Multimedia systems and equipment – Colour measurement and management – Part 1: Colour management – Default RGB colour space – sRGB.

IEC 61966-2-1 Amendment 1: ---1), Multimedia systems and equipment – Colour measurement and management – Part 2-1 Amendment 1: Colour management – Default RGB colour space – sRGB.

IEC 61966-3: 2000, Multimedia systems and equipment – Colour measurement and management – Part 3: Equipment using cathode ray tubes.

IEC 61966-4: 2000, Multimedia systems and equipment – Colour measurement and management – Part 4: Equipment using liquid crystal display panels.

IEC 61966-5: 2000, Multimedia systems and equipment – Colour measurement and management – Part 5: Equipment using plasma display panels.

Publication CIE 15.2: 1986, Colorimetry.

ITU-T Recommendation J.144, Objective perceptual video quality measurement techniques for digital cable television in the presence of a full reference.

ITU-R Recommendation BS.1387, Method for objective measurements of perceived audio quality.

ITU-T Recommendation P.931, Multimedia communications delays, synchronization and frame rate measurement.

---

1) Under development by TC 100/TA 2.  
2WD 2001-09-14

### 3 Terms and definitions

To understand this Technical Report, following terms and definitions apply.

#### 3.1

##### **audio-video communication system**

a system that handles audio, video and optionally other data streams in a synchronized way within users' perception in order to transmit and/or exchange information, which is assumed to operate over a local- or wide-area digital network

#### 3.2

##### **virtual conference**

meeting of a group of people who do not assemble to the same geographical place, but they exchange their views and opinions in use of multimedia logically connected each other

#### 3.3

##### **latency**

time required to send and receive a signal

#### 3.4

##### **linearity**

the number of video frames skipped at receiving end

#### 3.5

##### **PSNR**

peak-signal to noise ratio

#### 3.6

##### **PEAQ**

Perceived evaluation of audio quality

### 4 Configuration for quality assessment

#### 4.1 Input and output channels

Audio signal and video signal in audio-video streams shall be captured at the input and at the output channel, respectively, of the audio-video communication system as shown in figure 1.

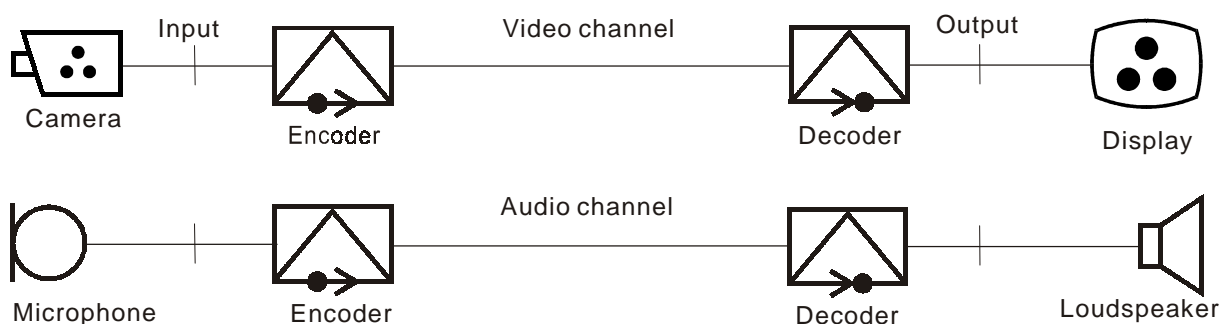


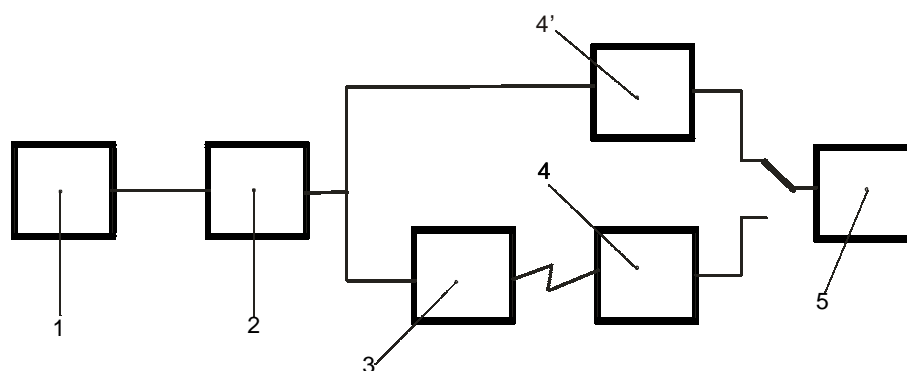
Figure 1 – Model of audio-video communication systems

#### 4.2 Points of input and output terminals

In the spirit of the end-to-end quality assessment of audio-video communication systems, the points for acquisition of raw data should be as far as ultimate end points as possible. However, since the methods of measurement and characterisation for equipment which incorporates

input transducers such as video cameras and microphones have already been standardised, such as in IEC 61146-1 and IEC 61146-2, and the methods of measurement and characterisation of equipment which incorporates output transducers such as video signal displays and loudspeakers, such as in IEC 61966-3, IEC 61966-4 and IEC 61966-5, they can be outside of the scope of the range of the end-to-end.

Figure 2 shows a schematic diagram for quality assessment under double-ended and full reference conditions.



#### Key

- 1 Original audio or video reference.
- 2 Preconditioner: Reduced dynamic range, frequency range for audio; reduced frame size and frame rate for video to fit to the quality assessment of the audio-video communication systems, if necessary.
- 3 Encoder for network streaming with a specified bit rate in order to fit to the bandwidth of end-to-end network connection.
- 4 Decoder and rendering for the received data to make them audible and visible.
- 4' Rendering for the preconditioned data to make them audible and visible, optional.
- 5 Data acquisition and calculation for quality assessment to provide information specified in this report.

**Figure 2 – Schematic diagram for quality assessment**

## 5 Video quality

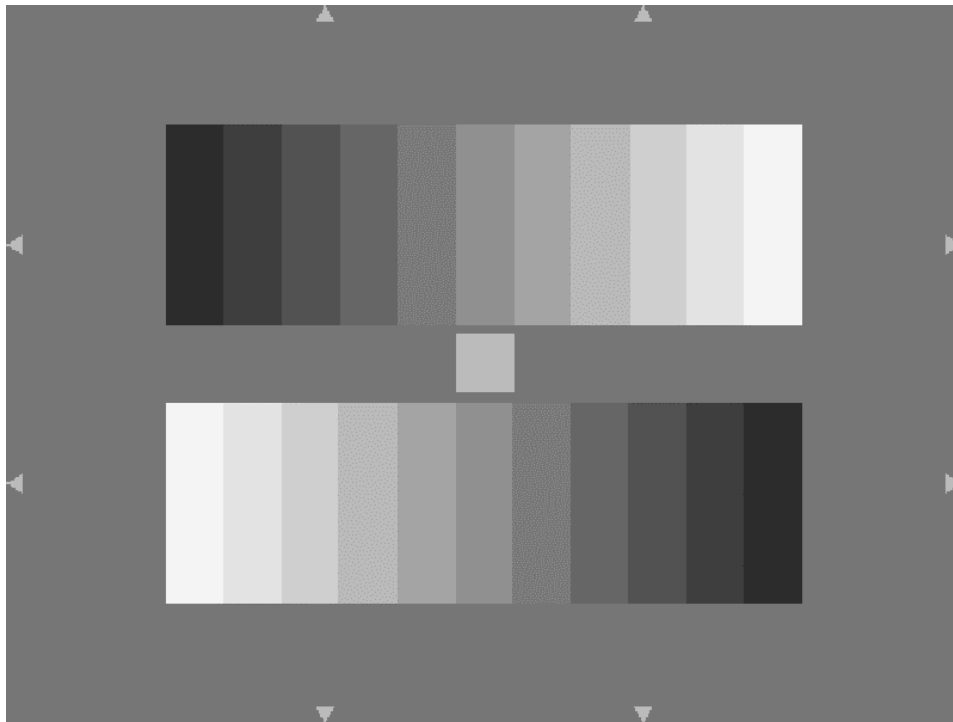
### 5.1 End-to-end tone reproduction

#### 5.1.1 Item to be assessed

End-to-end non-linearity in term of tone reproduction.

#### 5.1.2 Method of assessment

An image of the grey steps chart defined in IEC 61146-1, as shown in figure 3, should be used as the reference source at the item 1 in figure 2. The still neutral image should be prepared as a file for the item 2 in figure 2 and repeatedly encoded to be a streaming video transmitted to a network.



**Figure 3 – The image of the grey steps defined in IEC 61146-1**

Received streaming video should be decoded and rendered by a viewer for the incoming streaming videos. An image data to be displayed should be captured at an output terminal.

The image data should be compared in terms of three component data, R, G and B.

### 5.1.3 Form of reporting assessment result

The data for display versus the input image data should be reported as a table and a plot as shown in table 1 and figure 4, respectively, as examples, together with the audio-video communication system under assessment and the specification of the input-output point.

**Table 1 – An example of tone reproduction**

	Input image (%)	Output image (%)		
	$R_i = G_i = B_i$	$R_o$	$G_o$	$B_o$
0	2,0			
1	4,5			
2	8,1			
3	13,0			
4	19,8			
5	27,9			
6	37,8			
7	48,6			
8	63,0			
9	77,3			
10	89,9			



[Example plots should be filled here.]

**Figure 4 – An example plot of tone reproduction**

## **5.2 End-to-end colour reproduction**

### **5.2.1 Item to be assessed**

End-to-end colour shifts in the CIELAB colour space for static colour image.

### **5.2.2 Method of assessment**

An image of the colour reproduction chart defined in IEC 61146-1, as shown in figure 5, should be used as the reference source at the item 1 in figure 2. The still colour image should be prepared as a file for the item 2 in figure 2 and repeatedly encoded to be a streaming video transmitted to a network.



**Figure 5 – The image of the colour reproduction chart defined in IEC 61146-1**

Received streaming video should be decoded and rendered by a viewer for streaming videos. A colour image data to be displayed should be captured at an output terminal.

The image data should be acquired in terms of three component data, R, G and B.

### 5.2.3 Form of reporting assessment result

Input colours and output colours in R, G and B data should be regarded to be in sRGB defined in IEC 61966-2-1. They should be converted to CIE 1976 LAB uniform colour space. Colour differences  $\Delta E_{ab}^*$  between the reference data and the received data should be calculated and reported as shown in table 2 as an example.

**Table 2 – An example of colour reproduction**

	Input image			Output image			
	R <sub>i</sub> (%)	G <sub>i</sub> (%)	B <sub>i</sub> (%)	R <sub>o</sub> (%)	G <sub>o</sub> (%)	B <sub>o</sub> (%)	$\Delta E_{ab}^*$
0	87,053	80,546	87,216				
1	48,904	24,181	23,419				
2	37,405	27,352	12,466				
3	25,874	32,782	5,646				
4	12,176	34,717	19,279				
5	15,414	34,081	41,443				
6	17,982	29,222	61,449				
7	36,893	24,007	52,231				
8	51,332	22,896	45,507				
9	43,311	3,062	4,885				
10	83,988	56,759	4,964				
11	2,426	25,943	13,965				
12	3,259	7,178	18,424				
13	82,033	49,052	37,190				
14	10,356	12,908	4,612				

[Example plots should be filled here.]

**Figure 6 – An example plot of colour reproduction**

### 5.3 Peak-signal to noise ratio (PSNR)

#### 5.3.1 Item to be assessed

Peak-signal to noise power ratios, PSNR's, in three-dimensional coordinate systems.

#### 5.3.2 Method of assessment

##### 5.3.2.1 Definition of PSNR's in three-dimensional spaces

The peak-signal to noise ratio between a full reference image and a reproduced image recommended in ITU-T Recommendation J.144 should be used. It defined the PSNR by the following equation.

$$PSNR = 10 \log_{10} \left( \frac{S_{\max}^2}{MSE} \right) \quad (1)$$

$$MSE = \frac{1}{K} \sum_{p=P_1}^{P_2} \sum_{m=M_1}^{M_2} \sum_{n=N_1}^{N_2} (d(p, m, n) - o(p, m, n))^2$$

where  $K = \frac{1}{(P_2 - P_1 + 1)(M_2 - M_1 + 1)(N_2 - N_1 + 1)}$  ;  $d(p, m, n)$  and  $o(p, m, n)$  represent, respectively, degraded and original pixel vectors at frame  $p$ , row  $m$  and column  $n$ , and  $S_{\max}$  is the maximum possible value of the pixel vectors.

NOTE - PSNR requires a very high degree of normalisation to be used with confidence. The normalisation requires both spatial and temporal alignment as well as corrections for gain and offset.

For colour images, each picture element is normally composed of three dimensional values, red (R), green (G) and blue (B). Thus, following definition applies for the mean-square errors.

$$MSE_{RGB} = \frac{1}{K} \sum_{p=P_1}^{P_2} \sum_{m=M_1}^{M_2} \sum_{n=N_1}^{N_2} \left( (R_d - R_o)^2 + (G_d - G_o)^2 + (B_d - B_o)^2 \right) \quad (2)$$

where  $S_{\max(GB)} = 3 \times 2^{2(N-1)}$  for the values in  $N$ -bit encoding.

It is recommended to evaluate the PSNR in the more uniform colour space, CIE 1976 LAB, as follows.

$$MSE_{Lab} = \frac{1}{K} \sum_{p=P_1}^{P_2} \sum_{m=M_1}^{M_2} \sum_{n=N_1}^{N_2} \left( (L_d^* - L_o^*)^2 + (a_d^* - a_o^*)^2 + (b_d^* - b_o^*)^2 \right) \quad (3)$$

where  $S_{\max(Lab)} = \sqrt{(L_{\max}^*)^2 + (a_{\max}^*)^2 + (b_{\max}^*)^2}$ , actual value of which depends on a colour gamut of original RGB colour space. It is recommended to use the default RGB colour space defined by IEC 61966-2-1, in which  $S_{\max(Lab)} = 148,254$ .

It should be noted that the terms for summation in (3) are the square of the colour differences in the psychophysical uniform colour space defined in CIE 15.2. The average of colour differences  $\Delta E_{ab}^*$  should also be a metric of video quality. Namely,

$$\overline{\Delta E_{ab}^*} = \frac{1}{K} \sum_{p=P_1}^{P_2} \sum_{m=M_1}^{M_2} \sum_{n=N_1}^{N_2} \left( (L_d^* - L_o^*)^2 + (a_d^* - a_o^*)^2 + (b_d^* - b_o^*)^2 \right)^{\frac{1}{2}} \quad (4)$$

Additionally, luminance signal  $Y$  and two colour difference signals  $C_b$  and  $C_r$  denoted as  $Y_{cc}$  will also be calculated for comparison.

$$MSE_{Y_{cc}} = \frac{1}{K} \sum_{p=P_1}^{P_2} \sum_{m=M_1}^{M_2} \sum_{n=N_1}^{N_2} \left( (Y_d - Y_o)^2 + (C_{bd} - C_{bo})^2 + (C_{rd} - C_{ro})^2 \right) \quad (5)$$

where  $S_{\max}(Y_{cc}) = 1,01659$  in YCbCr system defined in IEC 61966-2-1 Amendment 1.

### 5.3.2.2 Configuration for quality assessment

It is recommended to make use of the commonly available video source as references such as the test sequences in the CRC as the original video reference for the item 1 in figure 2. Because of its high bit-rate and large frame size, the reference source should be reduced in frame size and bit-rate for use as the item 2 in figure 2, if necessary, for actual encoding as streaming video to a network with limited bandwidth.

It is needed to embed frame numbers at the item 2 in figure 2 so that they can be used to identify received frame numbers.

Since the values of PSNR's and additional metrics largely depend on video contents, varieties of video sources should be used for objective video quality assessment.

### 5.3.3 Form of reporting assessment results

The PSNR's in three-dimensional spaces Lab, Ycc and RGB together with the average colour difference  $\overline{\Delta E_{ab}^*}$  and one-dimensional PSNR's in  $L^*$  and  $Y$  should be reported as shown in figure 7.

The conditions of measurement such as frame size in pixels, frame rate, streaming bit-rate should also be reported.

NOTE – In order to demonstrate software developed by Chiba University in collaboration with Mitsubishi Electric Corp. for various quality metrics regarding the known hypothetical deterioration used in the Video Quality Expert Group (VQEG) in terms of three-dimensional PSNR's and one-dimensional PSNR's together with the average colour difference are attached in Annex A for information.

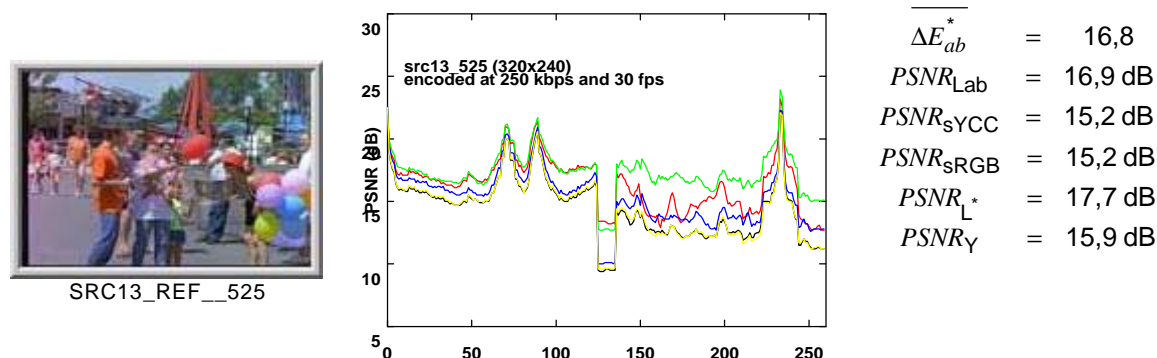


Figure 7a – An example of video quality assessment using the CRC SRC13\_REF\_\_525

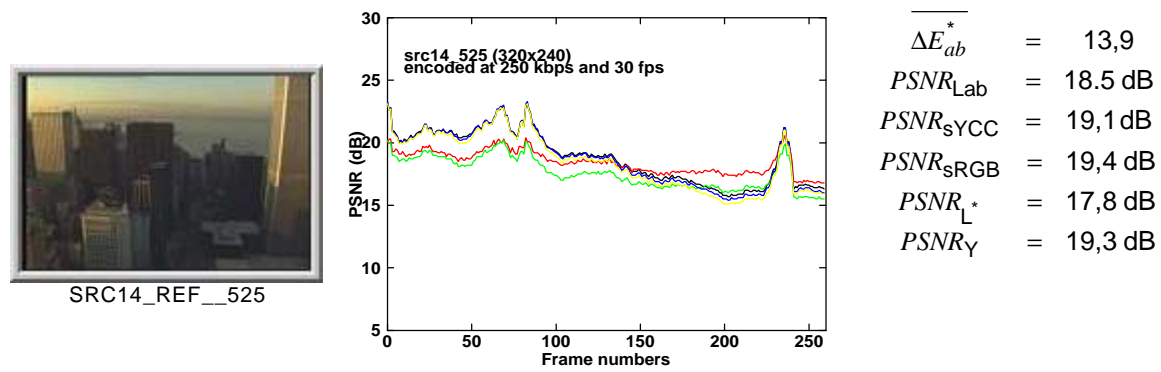


Figure 7b – An example of video quality assessment using the CRC SRC14\_REF\_\_525

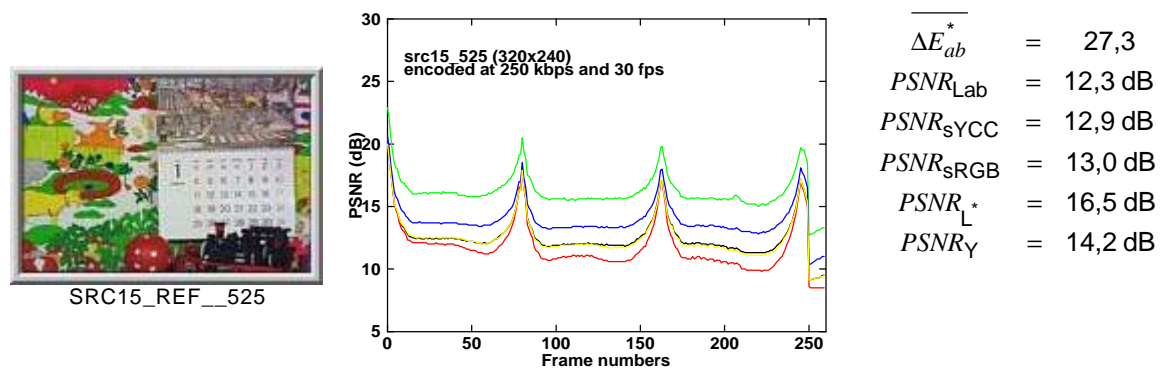


Figure 7c – An example of video quality assessment using the CRC SRC15\_REF\_\_525

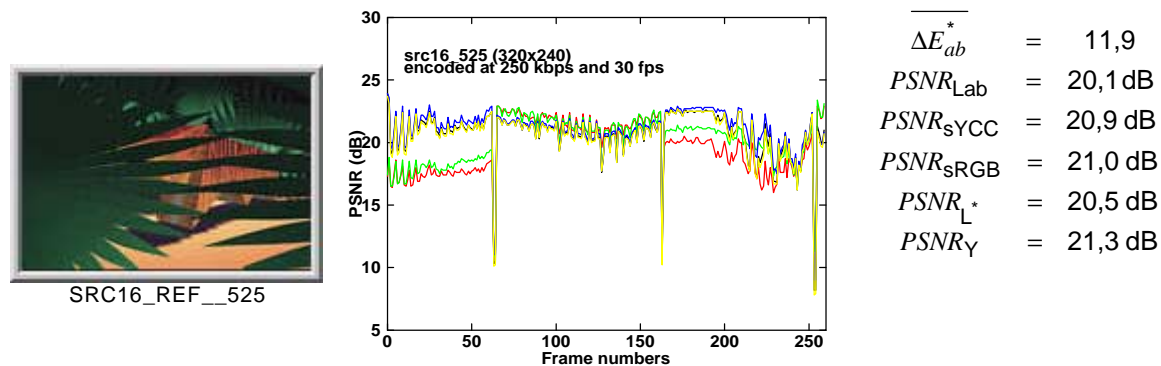


Figure 7d – An example of video quality assessment using the CRC SRC16\_REF\_\_525

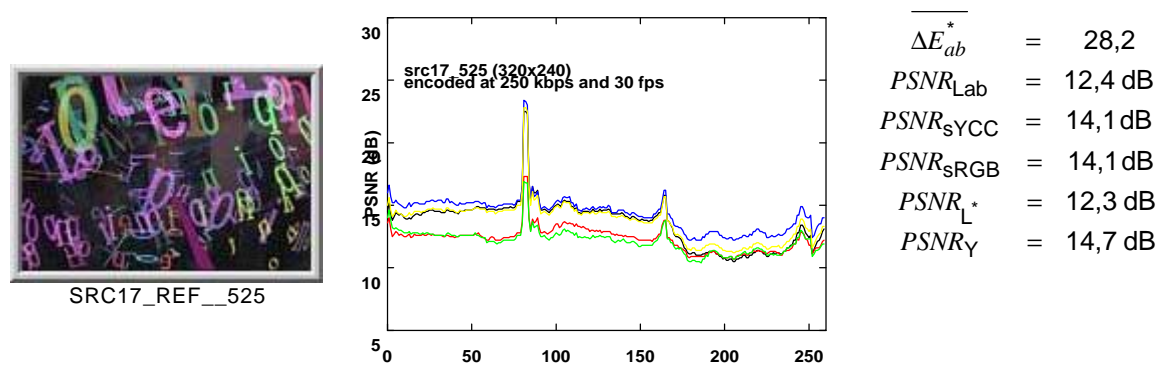


Figure 7e – An example of video quality assessment using the CRC SRC17\_REF\_\_525

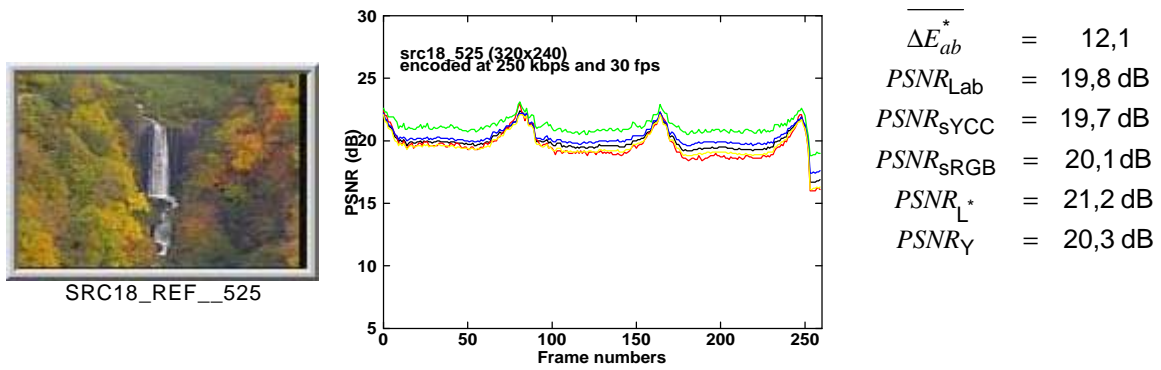


Figure 7f – An example of video quality assessment using the CRC SRC18\_REF\_\_525

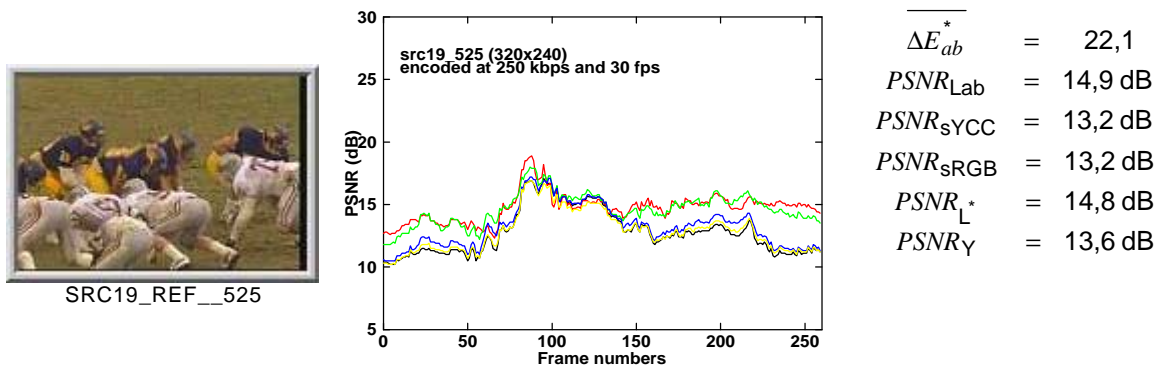


Figure 7g – An example of video quality assessment using the CRC SRC19\_REF\_\_525

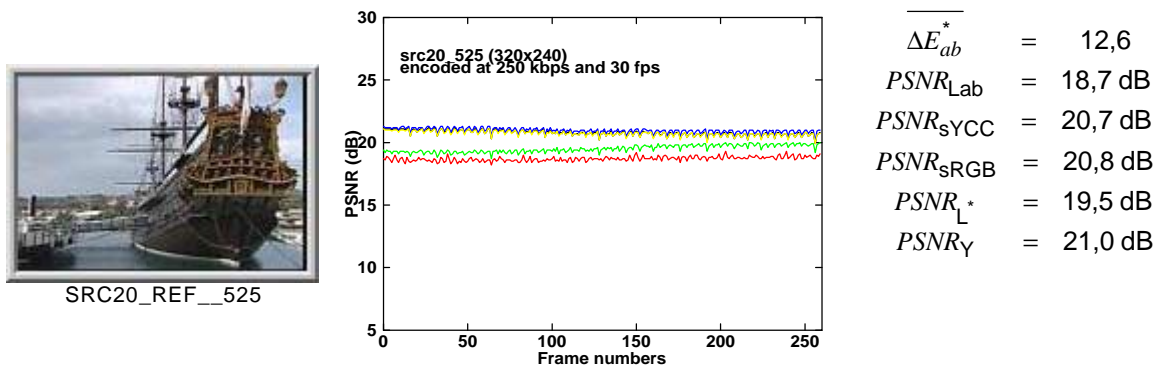


Figure 7h – An example of video quality assessment using the CRC SRC20\_REF\_\_525

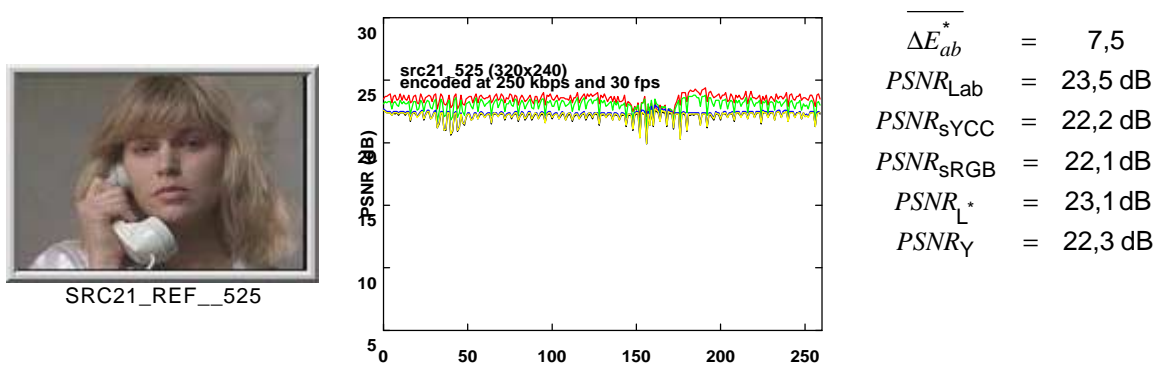
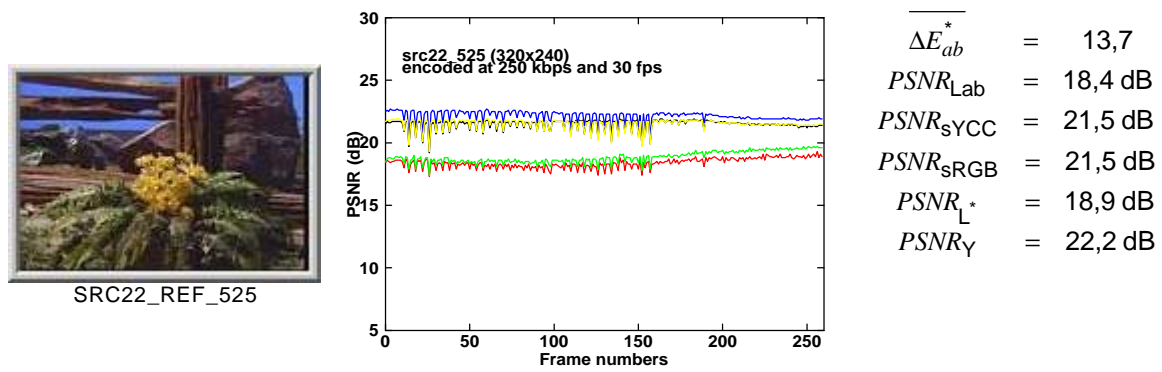


Figure 7i – An example of video quality assessment using the CRC SRC21\_REF\_\_525



**Figure 7j – An example of video quality assessment using the CRC SRC22\_REF\_\_525**

Condition of assessment:

- Video size: 320 x 240 pixels
- Frame rate: 30 fps
- Streaming bit rate: 250 kbps
- Network bandwidth: more than 250 kbps
- Reproduction: Microsoft Media Player® version 7.1

## 6 Audio quality

### 6.1 Perceived audio quality with full-reference signals

#### 6.1.1 Item to be assessed

Perceived evaluation of audio quality (PEAQ) recommended by ITU-R Recommendation BS.1387.

#### 6.1.2 Justification

Perceived audio quality (PEAQ) is one of the key factors when designing digital audio-video communication systems. Formal listening tests have been the relevant method for judging audio quality. However, subjective quality assessments are both time consuming and expensive. It was desirable to develop an objective measurement method in order to produce an estimate of the audio quality. Traditional objective measurement methods, like signal-to-noise-ratio (SNR) or total-harmonic-distortion (THD) have never really been shown to relate reliably to the perceived audio quality. The problems become even more evident when the methods are applied modern codecs, which are both non-linear and non-stationary. After through verification, ITU-R recommends an objective measurement method, known as PEAQ (Perceived Evaluation of Audio Quality), to estimate the perceived audio quality of equipment under test, for example a low bit-rate codec. This method is specified in ITU-R Recommendation BS.1387 and described briefly in Annex B.

The output variable from the PEAQ objective measurement method is the objective difference grade (ODG) and distortion index (DI). The ODG corresponds to the subjective difference grade (SDG) in the subjective domain. The resolution of the ODG is limited to one decimal. One should however be cautious and not generally expect that a difference between any pair of ODGs of a tenth of a grade is significant. The DI has the same meaning as the ODG. However, DI and ODG can only be compared quantitatively but not qualitatively. As a general rule, the ODG should be used as the quality measure for ODG values greater than approximately  $-3,6$ . The ODG correlates very well with subjective assessment in this range. When ODG value is less than  $-3,6$ , the DI should be used. Therefore, both ODG or DI variables shall be measured.

#### 6.1.3 Method of assessment and algorithm of PEAQ

The basic concept for PEAQ objective measurement method is illustrated in figure 8. It consists of two inputs, one for the (unprocessed) reference audio signal, corresponding to the item 2 of figure 2, and the other for the audio signal under the test. The latter may, for example, be the output signal of digital audio-video communication systems, corresponding to the output of the item 4 in figure 2, that is stimulated by the reference signal.

This measurement method is applicable to most types of audio signal processing equipment, both digital and analogue. It is, however, applied by focusing on digital audio communication channels in this document. The block "device under test" corresponds to the items 2 and 3 in figure 2.

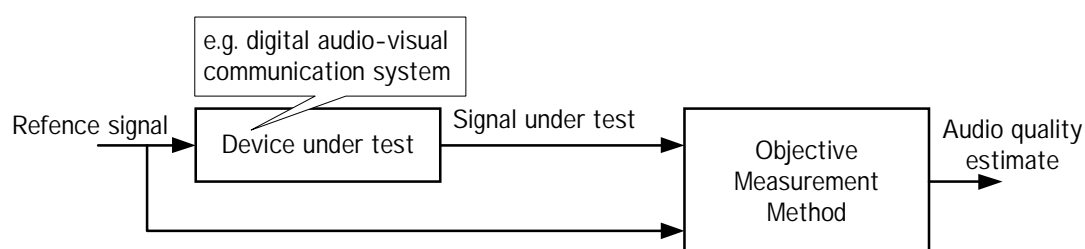


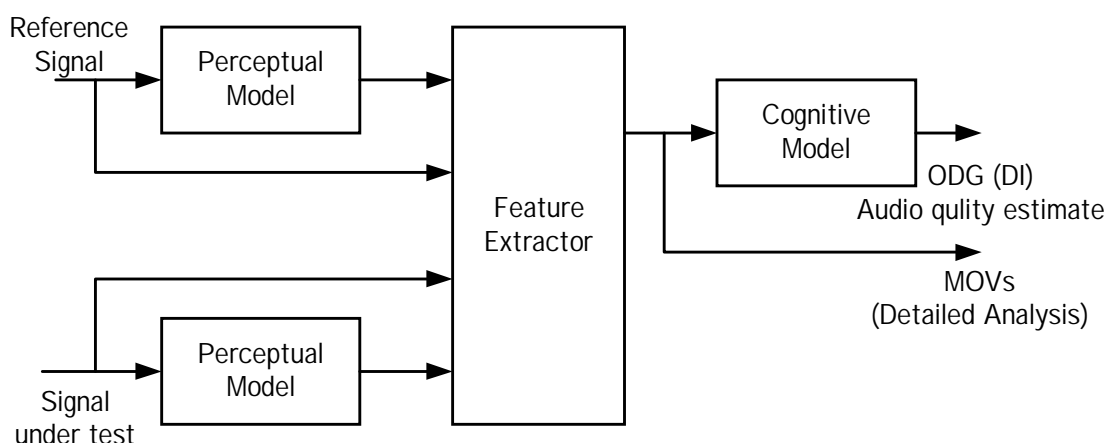
Figure 8 – Basic concept for making objective measurements



A representation of PEAQ model is shown in figure 9. The PEAQ method is based on generally accepted psychoacoustic principles. In general, it compares a signal that has been processed in some way with the corresponding time-aligned reference signal. In the first step of signal processing, the peripheral ear is modelled as known as “perceptual model”, or “ear model.” Concurrent frames of the reference and the processed signals are each transformed to the outputs of ear models. In a consecutive step, algorithm models the audible distortion present in the signal under test by comparing the outputs of the ear models. The information obtained by these processes results into several values, so called MOVs (model output variables) and may be useful for detailed analysis of the signal.

The final goal is to drive a quality metric, consisting of a single number that indicates the audibility of the distortions present in the signal under test. In order to archive this, some further processing of the MOVs is required which simulates the cognitive part of the human auditory system. Therefore, the PEAQ algorithm incorporates an artificial neural network.

There are two versions of the PEAQ; a “Basic” version, featuring a low complexity approach, and an “Advanced” version for higher accuracy at the trade off of higher complexity. The structure of both versions is very similar, and fits exactly into the PEAQ model shown in figure 9. The major difference between the basic and the advanced version is hidden in the respective ear models and the set of MOVs used. Annex B provides more information about PEAQ, which help readers to understand the measurement results.



**Figure 9 – Representation of PEAQ model**

It is recommended to use the reference items available from the ITU as WAV-files (Microsoft RIFF format) on a CD-ROM. All reference items have been sampled at 48 kHz in 16-bit PCM. The reference and test signals provided by the ITU have already been time and level adapted to each other, so that no additional gain or delay compensation is required.

The measurement algorithm must be adjusted to a listening level of 92 dB SPL.

#### 6.1.4 Form of reporting assessment results

PEAQ measurement results should be reported in terms of the names of the reference and signal under test items<sup>2</sup>, and the resulting DI and ODG values in a table.

Table 3 is related to the basic version, and table 4 contains the values for the advanced version.

<sup>2</sup> The names of the corresponding reference items are derived by replacing the substring “cod” in the names of the test items by “ref”, e.g. the reference item for “bcodtri.wav” is “breftri.wav”.

**Table 3 – Test items and resulting DI and ODG values for the basic version**

Item	DI	ODG	Item	DI	ODG	Item	DI	ODG
Acodsna.wav	1,304	-0,7	Fcodtr2.wav	-0,045	-1,9	lcodhrp.wav	1,041	-0,9
Bcodtri.wav	1,949	-0,3	fcodtr3.wav	-0,715	-2,6	lcodpip.wav	1,973	-0,3
Ccodsax.wav	0,048	-1,8	gcodcla.wav	1,781	-0,4	mcodcla.wav	-0,436	-2,3
Dcodryc.wav	1,648	-0,5	hcodryc.wav	2,291	-0,2	ncodsfe.wav	3,135	0,0
Ecodsmg.wav	1,731	-0,4	Hcodstr.wav	2,403	-0,1	scodclv.wav	1,689	-0,4
Fcodsb1.wav	0,677	-1,2	lcodsna.wav	-3,029	-3,8			
Fcodtr1.wav	1,419	-0,6	kcodsme.wav	3,093	0,0			

**Table 4 – Test items and resulting DI and ODG values for the advanced version**

Item	DI	ODG	Item	DI	ODG	Item	DI	ODG
Acodsna.wav	1,632	-0,5	Fcodtr2.wav	0,162	-1,7	Lcodhrp.wav	1,538	-0,5
Bcodtri.wav	2,000	-0,3	Fcodtr3.wav	-0,783	-2,7	Lcodpip.wav	2,149	-0,2
Ccodsax.wav	0,567	-1,3	Gcodcla.wav	1,457	-0,6	Mcodcla.wav	0,430	-1,4
Dcodryc.wav	1,725	-0,4	Hcodryc.wav	2,410	-0,1	Ncodsfe.wav	3,163	0,0
Ecodsmg.wav	1,594	-0,5	Hcodstr.wav	2,232	-0,2	Scodclv.wav	1,972	-0,3
Fcodsb1.wav	1,039	-0,9	lcodsna.wav	-2,510	-3,7			
Fcodtr1.wav	1,555	-0,5	Kcodsme.wav	2,765	-0,0			

## 6.2 Sampling rate and quantization resolution

### 6.2.1 Item to be assessed

Sampling rate and bandwidth of the reference and the processed audio signals.

### 6.2.2 Method of assessment

Sampling rate is relevant to the bandwidth of audio signals. For high-quality audio signals, the sampling rates 48 kHz, 44,1 kHz are used. The sampling rate and bandwidth of the reference and processed audio signals should be extracted.

Resolution of quantization relates to the dynamic range of audio signals or quantization noise. For high-quality audio signals, the linear (or uniform) quantization method, which have 16-bit quantization resolution, are used. The resolution and quantization method should be identified.

### 6.2.3 Form of reporting assessment results

Extracted and identified values should be reported.

### **6.3 Delay**

#### **6.3.1 Item to be assessed**

Delay time in seconds from audio inputs to an encoder and received digital audio signals.

#### **6.3.2 Method of assessment**

Pulsed audio signals should be used as input in terms of the item 2 of figure 2. Lap time between the input of the item 3 and the output of the item 4 in figure 2 should be measured in seconds.

NOTE – In most audio communication systems over the digital networks, a buffering scheme is incorporated. Therefore, buffering time is also taken into account in the measurement.

#### **6.3.3 Form of reporting assessment result**

The measured delay time should be reported in seconds.

## 7 Total quality

### 7.1 Synchronization of audio and video (lip sync)

#### 7.1.1 Item to be assessed

Temporal synchronization between audio and video channels.

#### 7.1.2 Method of assessment

The raison d'être of true multimedia systems, in contrast to a mere collection of unrelated media channels, is the ability to keep temporal synchronization among different channels. It is therefore vitally important to include a temporal synchronization quality measurement in the quality assessment items of audio-video communication systems.

The framework for measurement of temporal synchronization among media channels is given in ITU-T Recommendation P.931. Its basic assumption is that media signal can be captured at such interfaces as the camera output and the display input for the visual channel and the microphone output and the loudspeaker input for the audio channel. This assumption is shown in figure 1.

Media signals at those interfaces are digitised, if necessary, broken into fixed sized frames, and given timestamps. For the details for this procedure, refer to ITU-T Recommendation P.931.

The audio and the video media streams being considered, digitized frames of these media streams are given sequence numbers as follows:

- $A(m)$  and  $V(n)$  are the input audio and the video frames, respectively ( $m$  and  $n$  are the sequence numbers for each stream). It is assumed that they are associated in the sense that they correspond to the same event.
- $A'(p)$  and  $V'(q)$  are the output audio and the video frames, respectively.
- $T_A(m)$  and  $T'_A(n)$  are the timestamps for  $A(m)$  and  $A'(n)$ , respectively. Timestamps for other frames are defined in the same way.

For each input frame, however not all the input frames are to be used as described in ITU-T Recommendation P.931, the corresponding output frame is to be found. Since the media stream data is modified, distorted, dropped and reshaped, the matching process is non-trivial. For video frames, the method which utilizes the metrics of PSNR's assessed in 5 will be used. For audio frames, a two-stage process, which employs the audio envelopes for a coarse stage and the power spectral densities for a fine stage, will be used. For details, refer to ITU-T Recommendation P.931.

It is assumed here that the matching relationships between  $A(m)$  and  $A'(p)$ , and  $V'(n)$  and  $V'(q)$  are established. Under this assumption, the time skew between the audio and video frames is given by the following expression:

$$S_{AV}(p, q) = O'_{AV}(p, q) - O_{AV}(m, n) \quad (6)$$

where  $Q_{AV}(m, n) = T_A(m) - T_V(n)$  and  $Q'_{AV}(p, q) = T'_A(p) - T'_V(q)$ .

NOTE 1 – It is important to choose appropriate input audio signals to obtain a valid and meaningful assessment result. If the video signal contains still or almost-still scenes, the process for matching the input and the output frames will become difficult or even impossible. A similar caution is to be applied to the assessment of the audio channel.

NOTE 2 – Modern video compression schemes give fluctuating compression, transmission (when variable bit rate encoding is employed) and decompression times, depending on the input properties. Therefore, the appropriate input signals suited for the assumed application should be used for assessment.

NOTE 3 – For systems with a low video frame rate, it is sometimes natural and desirable to have a larger temporal skew between the video and the audio streams, because the video delay fluctuates while the audio data generally flows isochronously.

Selection of standard audio-video input streams suited for common usage is left for future study.

### 7.1.3 Form of reporting assessment results

The report from the measurement should be expressed such that any variation between individual measurement is clearly illustrated. Classical summary statistics (for example, minimum, maximum, average and standard deviation) may also be reported.

## 7.2 Scalability

### 7.2.1 Item to be assessed

Autonomous function to tune frame rate dynamically depending on available bandwidth between the sender and the receiver.

### 7.2.2 Method of assessment

The method for measurement of scalability is left for future study.

### 7.2.3 Form of reporting assessment results

Under consideration.

## 7.3 Overall quality

### 7.3.1 Item to be assessed

Overall quality factor in terms of audio and video interaction.

### 7.3.2 Method of assessment

Overall quality of audio-video communication systems  $OQ_{AV}$  should be defined as follows.

$$OQ_{AV} = aQ_V + bQ_A + cQ_{V\&A} \quad (7)$$

where  $Q_V$  is the objective quality metric assessed in 5,  $Q_A$  is the objective quality metric assessed in 6,  $Q_{V\&A}$  is the objective quality metric assessed in 7; the coefficients  $a$ ,  $b$  and  $c$  are the weighting factors, which depend on actual applications of the audio-video communication system.

### 7.3.3 Form of reporting assessment results

The overall quality factor should be reported with sufficient information on the audio-video communication system under assessment.

## **Annex A (informative)**

### **PSNR's defined in three-dimensional spaces applied to hypothetical deterioration over the reference video sources**

#### **A.1 Introduction**

This informative annex is intended to demonstrate the definitions PSNR's in three-dimensional vector space for each of pixel that consists of a frame of videos. The definition for PSNR in the CIELAB is given in equation (3), PSNR in sYCC in equation (5), PSNR in sRGB in equation (2). The average colour difference defined in equation (4) is also included in this annex for comparison together with one-dimensional PSNR's in  $L^*$  and  $Y$ .

The values of the objective quality measures will be easily compared with other possible future measures and the results of subjective assessment of video quality.

#### **A.2 Test sources and hypothetical deterioration**

In this annex, known 16 different hypothetical deterioration over the digital video files in prepared in ITU-R BT.601-5 format and used in the Video Quality Expert Group (VQEG). The source videos are labelled from SRC13\_REF\_\_525.yuv to SRC22\_REF\_\_525. They are made use of by the permission of the VQEG.

Software for varieties of objective measures have been developed in Chiba University, Japan, in collaboration with Mitsubishi Electric Corp. The values have been obtained for reduced frame size of 320 x 240 pixels per frame over 260 frames. In other words,  $P1=1$ ,  $P2=260$ ,  $M1=1$ ,  $M2=240$  and  $N1=1$ ,  $N2=320$  in equation (1). Numerical results are shown in tables A.1 to A.5.

**Table A.1 – PSNR's in various colour spaces and the colour difference for SRC13 and SRC14**

	Lab	sYCC	sRGB	L*	Y	$\Delta E$		Lab	SYCC	sRGB	L*	Y	$\Delta E$
<b>hrc1/src13</b>	20.5	23.2	23.6	26.3	26.3	8.3	<b>hrc1/src14</b>	22.4	25.8	25.9	26.6	28.1	7.5
<b>hrc2/src13</b>	23.6	23.5	23.2	25.9	25.0	5.4	<b>hrc2/src14</b>	25.7	24.3	24.5	25.4	24.3	4.9
<b>hrc3/src13</b>	22.2	22.7	22.3	25.6	24.6	5.8	<b>hrc3/src14</b>	24.9	23.8	24.0	25.1	24.1	4.7
<b>hrc4/src13</b>	21.4	22.1	21.7	25.6	24.7	7.4	<b>hrc4/src14</b>	24.6	23.9	24.0	25.4	24.3	5.5
<b>hrc5/src13</b>	20.4	19.3	19.0	21.2	20.3	8.0	<b>hrc5/src14</b>	22.5	19.7	20.0	20.7	19.6	5.9
<b>hrc6/src13</b>	22.2	22.6	22.1	25.9	24.9	6.2	<b>hrc6/src14</b>	24.5	23.6	23.8	25.3	24.1	5.0
<b>hrc7/src13</b>	22.2	21.1	20.7	23.0	22.1	5.9	<b>hrc7/src14</b>	24.5	21.5	21.7	22.5	21.4	4.1
<b>hrc8/src13</b>	21.9	22.3	21.9	25.3	24.5	6.7	<b>hrc8/src14</b>	24.3	23.5	23.7	24.9	24.0	5.3
<b>hrc9/src13</b>	21.6	20.6	20.3	22.8	21.8	6.9	<b>hrc9/src14</b>	24.3	21.4	21.7	22.4	21.4	4.5
<b>hrc10/src13</b>	22.1	20.9	20.6	23.0	22.0	6.3	<b>Hrc10/src14</b>	24.3	21.4	21.6	22.5	21.4	4.4
<b>hrc11/src13</b>	21.7	22.8	22.5	24.5	25.3	6.9	<b>Hrc11/src14</b>	25.5	26.0	26.1	24.6	26.3	4.1
<b>hrc12/src13</b>	22.4	23.6	23.3	24.8	26.0	5.9	<b>Hrc12/src14</b>	26.0	26.2	26.4	24.8	26.4	3.7
<b>hrc13/src13</b>	21.3	20.7	20.6	23.4	22.2	6.8	<b>Hrc13/src14</b>	21.5	20.8	21.7	23.2	21.7	5.4
<b>hrc14/src13</b>	21.2	20.3	20.0	22.7	21.6	7.9	<b>Hrc14/src14</b>	23.9	21.3	21.6	22.4	21.3	5.3
<b>hrc15/src13</b>	21.9	22.1	21.7	25.3	24.4	7.6	<b>Hrc15/src14</b>	25.8	25.8	26.0	27.2	26.3	5.6
<b>hrc16/src13</b>	22.1	22.8	22.3	25.8	25.2	7.0	<b>Hrc16/src14</b>	26.0	26.0	26.2	27.4	26.5	5.3

NOTE 1 – hrc16/src14 and so on correspond to hypothetically degraded video (hrc16) from the reference source video (src14), respectively.

NOTE 2 – All videos are in size of 320 x 240 pixels, each of which has 24-bit colour depth.

**Table A.2 – PSNR's in various colour spaces and the colour difference for SRC15 and SRC16**

	Lab	sYCC	sRGB	L*	Y	$\Delta E$		Lab	SYCC	sRGB	L*	Y	$\Delta E$
<b>hrc1/src15</b>	11.8	13.6	13.1	20.7	19.5	24.8	<b>hrc1/src16</b>	20.3	21.6	21.8	23.8	25.7	9.5
<b>hrc2/src15</b>	17.0	18.5	18.4	24.2	23.1	10.8	<b>hrc2/src16</b>	27.1	28.1	28.0	31.1	32.0	4.4
<b>hrc3/src15</b>	15.1	16.7	16.5	23.1	21.7	13.2	<b>hrc3/src16</b>	29.2	29.0	28.9	31.0	31.9	2.4
<b>hrc4/src15</b>	13.6	15.0	14.5	23.0	21.2	18.7	<b>hrc4/src16</b>	22.9	23.7	23.6	28.3	28.3	6.0
<b>hrc5/src15</b>	14.6	15.4	15.2	19.8	18.9	15.9	<b>hrc5/src16</b>	21.7	22.0	21.9	24.8	25.5	6.0
<b>hrc6/src15</b>	14.0	15.4	15.0	23.0	21.3	17.3	<b>hrc6/src16</b>	23.5	24.2	24.0	28.7	28.7	5.1
<b>hrc7/src15</b>	16.1	17.0	16.9	21.1	20.3	12.1	<b>hrc7/src16</b>	22.8	22.9	22.8	25.7	26.4	4.4
<b>hrc8/src15</b>	14.0	15.3	15.0	22.6	20.9	17.6	<b>hrc8/src16</b>	23.4	24.2	24.0	28.4	28.5	5.2
<b>hrc9/src15</b>	15.9	16.6	16.5	20.7	19.7	13.1	<b>hrc9/src16</b>	22.8	22.7	22.7	25.5	26.1	4.6
<b>hrc10/src15</b>	16.4	17.3	17.2	22.1	21.0	11.9	<b>hrc10/src16</b>	24.9	25.8	25.4	28.9	30.3	3.8
<b>hrc11/src15</b>	15.8	17.2	17.0	22.9	22.0	12.8	<b>hrc11/src16</b>	25.4	27.5	27.3	27.8	31.6	3.8
<b>hrc12/src15</b>	16.0	17.5	17.3	23.3	22.7	12.0	<b>hrc12/src16</b>	25.7	27.9	27.6	28.0	32.2	3.5
<b>hrc13/src15</b>	15.4	16.2	16.1	21.5	19.7	14.6	<b>hrc13/src16</b>	23.3	23.5	23.6	29.1	29.5	4.3
<b>hrc14/src15</b>	15.6	16.1	16.0	20.6	19.2	14.3	<b>hrc14/src16</b>	22.9	22.6	22.6	25.2	25.4	5.2
<b>hrc15/src15</b>	15.7	16.1	16.0	21.1	19.3	15.4	<b>hrc15/src16</b>	23.7	23.3	23.5	26.0	26.2	5.8
<b>hrc16/src15</b>	15.7	16.3	16.2	21.6	19.9	14.9	<b>hrc16/src16</b>	23.9	23.5	23.7	26.2	26.5	5.6

NOTE 1 – hrc16/src14 and so on correspond to hypothetically degraded video (hrc16) from the reference source video (src14), respectively.

NOTE 2 – All videos are in size of 320 x 240 pixels, each of which has 24-bit colour depth.

**Table A.3 – PSNR's in various colour spaces and the colour difference for SRC17 and SRC18**

	Lab	sYCC	sRGB	L*	Y	ΔE		Lab	sYCC	sRGB	L*	Y	ΔE
<b>hrc1/src17</b>	15.8	19.2	19.2	20.8	23.3	16.7	<b>hrc1/src18</b>	18.3	21.0	20.7	23.2	25.6	10.2
<b>hrc2/src17</b>	20.2	23.2	23.6	26.6	26.9	9.2	<b>hrc2/src18</b>	22.8	24.8	24.5	28.0	28.7	6.0
<b>hrc3/src17</b>	20.2	23.2	23.3	26.2	27.1	8.3	<b>hrc3/src18</b>	22.4	24.2	23.8	27.7	28.0	6.5
<b>hrc4/src17</b>	18.6	21.2	21.6	25.2	25.0	11.1	<b>hrc4/src18</b>	18.1	20.4	19.7	26.6	26.9	9.9
<b>hrc5/src17</b>	18.0	20.1	20.5	22.7	23.0	11.8	<b>hrc5/src18</b>	18.9	20.1	20.0	21.7	22.6	9.0
<b>hrc6/src17</b>	18.5	20.8	21.1	24.9	24.7	10.0	<b>hrc6/src18</b>	19.3	21.6	21.0	27.2	27.4	8.4
<b>hrc7/src17</b>	19.8	21.8	22.1	24.5	24.9	8.6	<b>hrc7/src18</b>	20.3	21.5	21.5	22.8	23.7	7.2
<b>hrc8/src17</b>	18.1	20.5	20.8	24.3	24.2	10.9	<b>hrc8/src18</b>	19.5	21.7	21.2	26.8	27.2	8.4
<b>hrc9/src17</b>	18.6	20.6	20.9	23.4	23.7	10.3	<b>hrc9/src18</b>	20.4	21.5	21.6	22.8	23.6	7.4
<b>hrc10/src17</b>	19.7	21.9	22.2	24.8	25.2	8.9	<b>hrc10/src18</b>	21.6	23.0	22.8	25.3	26.0	6.5
<b>hrc11/src17</b>	18.1	20.4	20.7	23.2	23.9	10.8	<b>hrc11/src18</b>	21.5	24.4	24.0	26.7	29.8	6.4
<b>hrc12/src17</b>	19.0	21.4	21.8	24.1	25.1	9.3	<b>hrc12/src18</b>	21.7	24.5	24.1	26.7	30.1	6.1
<b>hrc13/src17</b>	16.9	18.6	19.0	22.0	22.0	13.2	<b>hrc13/src18</b>	21.9	24.0	23.6	27.7	27.9	6.9
<b>hrc14/src17</b>	18.2	20.3	20.6	23.4	23.4	11.6	<b>hrc14/src18</b>	21.3	22.8	22.6	25.2	25.7	7.2
<b>hrc15/src17</b>	17.8	20.0	20.4	23.0	23.0	13.4	<b>hrc15/src18</b>	21.6	23.7	23.3	29.0	28.5	7.8
<b>hrc16/src17</b>	18.2	20.7	21.2	23.8	24.0	12.5	<b>hrc16/src18</b>	21.7	23.9	23.4	29.4	29.2	7.4

NOTE 1 – hrc16/src14 and so on correspond to hypothetically degraded video (hrc16) from the reference source video (src14), respectively.

NOTE 2 – All videos are in size of 320 x 240 pixels, each of which has 24-bit colour depth.

**Table A.4 – PSNR's in various colour spaces and the colour difference for SRC19 and SRC20**

	Lab	sYCC	sRGB	L*	Y	ΔE		Lab	sYCC	sRGB	L*	Y	ΔE
<b>hrc1/src19</b>	20.0	22.6	22.6	23.2	25.6	7.8	<b>hrc1/src20</b>	15.8	17.4	17.2	20.1	20.2	12.7
<b>hrc2/src19</b>	23.6	25.1	24.9	27.9	28.6	4.8	<b>hrc2/src20</b>	20.6	20.8	20.7	23.7	22.1	6.9
<b>hrc3/src19</b>	23.1	24.6	24.4	27.7	28.0	5.8	<b>hrc3/src20</b>	18.7	19.3	19.3	22.6	21.3	8.2
<b>hrc4/src19</b>	20.3	22.5	22.1	26.6	27.1	6.9	<b>hrc4/src20</b>	18.7	19.2	19.0	22.6	21.0	8.8
<b>hrc5/src19</b>	19.8	20.8	20.7	22.5	23.3	7.4	<b>hrc5/src20</b>	18.7	16.2	16.0	18.5	16.6	8.3
<b>hrc6/src19</b>	20.6	22.7	22.3	27.0	27.2	6.6	<b>hrc6/src20</b>	18.8	19.4	19.2	23.1	21.4	8.1
<b>hrc7/src19</b>	21.0	21.7	21.7	22.8	23.5	5.9	<b>hrc7/src20</b>	19.4	17.5	17.3	19.6	18.1	7.2
<b>hrc8/src19</b>	20.7	22.6	22.3	26.6	26.9	6.8	<b>hrc8/src20</b>	18.6	19.2	19.0	22.8	21.2	8.4
<b>hrc9/src19</b>	21.2	23.1	22.7	27.2	27.2	6.2	<b>hrc9/src20</b>	20.0	20.3	20.1	23.5	22.0	6.4
<b>hrc10/src19</b>	20.0	21.3	21.1	24.4	24.7	7.8	<b>Hrc10/src20</b>	20.3	18.8	18.6	21.3	19.5	6.5
<b>hrc11/src19</b>	21.4	23.6	23.3	25.6	27.7	6.3	<b>Hrc11/src20</b>	19.9	21.5	21.4	23.3	23.8	6.6
<b>hrc12/src19</b>	22.4	24.7	24.4	25.9	28.8	5.4	<b>Hrc12/src20</b>	20.3	21.9	21.7	23.4	24.1	6.2
<b>hrc13/src19</b>	20.8	21.8	21.7	24.2	24.3	6.9	<b>Hrc13/src20</b>	19.8	18.8	18.7	21.4	19.8	7.8
<b>hrc14/src19</b>	21.1	22.1	22.0	24.8	25.1	7.1	<b>Hrc14/src20</b>	19.6	18.4	18.2	21.1	19.3	7.6
<b>hrc15/src19</b>	23.3	24.6	24.4	28.8	28.5	5.9	<b>Hrc15/src20</b>	19.5	20.3	20.2	23.5	22.2	8.4
<b>hrc16/src19</b>	23.6	25.1	24.8	29.4	29.4	5.4	<b>Hrc16/src20</b>	19.6	20.4	20.4	23.7	22.4	8.3

NOTE 1 – hrc16/src14 and so on correspond to hypothetically degraded video (hrc16) from the reference source video (src14), respectively.

NOTE 2 – All videos are in size of 320 x 240 pixels, each of which has 24-bit colour depth.



**Table A.5 – PSNR's in various colour spaces and the colour difference for SRC21 and SRC22**

	Lab	sYCC	sRGB	L*	Y	$\Delta E$		Lab	sYCC	sRGB	L*	Y	$\Delta E$
<b>Hrc1/src21</b>	23.1	25.3	25.8	22.8	25.8	5.9	<b>hrc1/src22</b>	14.6	18.0	17.6	22.3	24.1	16.8
<b>Hrc2/src21</b>	29.3	29.1	29.2	28.5	29.6	3.2	<b>hrc2/src22</b>	18.9	21.7	21.0	26.3	26.5	9.3
<b>Hrc3/src21</b>	29.4	28.8	28.8	28.4	29.3	2.9	<b>hrc3/src22</b>	17.0	19.9	19.3	24.8	24.9	11.0
<b>Hrc4/src21</b>	28.4	27.7	27.9	27.3	28.2	3.5	<b>hrc4/src22</b>	17.4	20.1	19.4	25.4	25.6	11.2
<b>Hrc5/src21</b>	25.7	24.0	24.1	22.8	24.0	3.5	<b>hrc5/src22</b>	17.4	18.9	18.0	21.5	21.9	11.2
<b>Hrc6/src21</b>	29.5	28.3	28.5	27.9	28.6	2.8	<b>hrc6/src22</b>	17.2	20.0	19.3	25.7	25.8	10.8
<b>Hrc7/src21</b>	26.0	24.4	24.5	23.1	24.4	3.0	<b>hrc7/src22</b>	18.0	19.9	19.2	22.8	23.3	9.8
<b>Hrc8/src21</b>	29.1	28.1	28.3	27.5	28.4	3.0	<b>hrc8/src22</b>	17.2	19.9	19.2	25.1	25.2	11.1
<b>Hrc9/src21</b>	30.7	29.4	29.5	28.5	29.6	2.0	<b>hrc9/src22</b>	17.9	20.5	19.8	25.4	25.5	9.7
<b>hrc10/src21</b>	28.5	26.9	27.0	25.8	26.9	2.5	<b>hrc10/src22</b>	18.2	20.3	19.5	23.9	24.2	9.7
<b>hrc11/src21</b>	28.8	30.6	30.7	26.7	31.0	2.4	<b>hrc11/src22</b>	18.0	20.8	20.3	24.4	25.6	10.0
<b>hrc12/src21</b>	28.9	30.8	30.9	26.7	31.2	2.2	<b>hrc12/src22</b>	18.3	21.3	20.7	24.9	26.5	9.3
<b>hrc13/src21</b>	27.4	25.8	25.9	25.0	25.9	3.2	<b>hrc13/src22</b>	16.9	18.9	18.5	22.7	22.6	12.2
<b>hrc14/src21</b>	28.2	26.7	26.8	25.7	26.8	2.9	<b>hrc14/src22</b>	17.8	19.7	19.0	23.2	23.2	11.0
<b>hrc15/src21</b>	30.5	30.4	30.5	30.3	31.1	3.2	<b>hrc15/src22</b>	17.8	20.2	19.8	24.4	23.9	12.0
<b>hrc16/src21</b>	30.6	30.5	30.6	30.4	31.2	3.2	<b>hrc16/src22</b>	18.1	20.8	20.3	25.4	25.1	11.3

NOTE 1 – hrc16/src14 and so on correspond to hypothetically degraded video (hrc16) from the reference source video (src14), respectively.

NOTE 2 – All videos are in size of 320 x 240 pixels, each of which has 24-bit colour depth.

## Annex B (informative)

### PEAQ objective measurement method outline

#### B.1 Basic concept of the PEAQ measurement algorithm

The basic concept for PEAQ objective measurement method is illustrated in figure B.1. It consists of two inputs, one for the (unprocessed) reference signal and one for the signal under the test. The latter may for example be the output signal of the codec that is stimulated by the reference signal.

This measurement method is applicable to most types of audio signal processing equipment, both digital and analogue. It is, however, expected that many applications will focus on audio codecs.

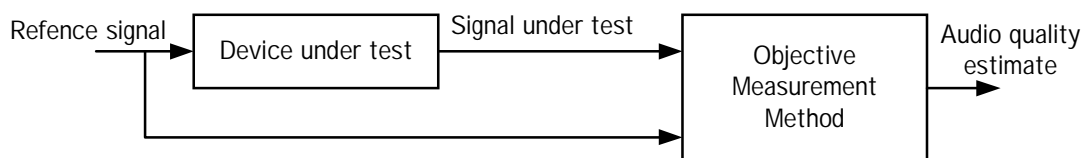


Figure B.1 – Basic concept for making objective measurements

A high-level representation of PEAQ model is shown in figure B.2. PEAQ method is based on generally accepted psychoacoustic principles. In general it compares a signal that has been processed in some way with the corresponding time-aligned reference signal. In the first signal processing step the peripheral ear is modelled (“perceptual model”, or “ear model”). Concurrent frames of the reference and processed signal are each transformed to the outputs of ear models. In a consecutive step, algorithm models the audible distortion present in the signal under test by comparing the outputs of the ear models. The information obtained by these process results into several values, so called MOVs (“Model Output Variables”) and may be useful for detailed analysis of the signal.

The final goal instead is to drive a quality measure, consisting of a single number that indicates the audibility of the distortions present in the signal under test. In order to archive this, some further processing of the MOVs is required which simulates the cognitive part of the human auditory system. Therefore the PEAQ algorithm uses an artificial neural network.

There are two versions of PEAQ, a “Basic” version, featuring a low complexity approach, and an “Advanced” version for higher accuracy at the trade off of higher complexity. The structure of both versions is very similar, and fits exactly into the PEAQ model shown in figure.B.2. The major differences between the “Basic” and the “Advanced” version are hidden in the respective ear models and the set of MOVs used. The “Basic” and “Advanced” versions are described in B.2 and B.3

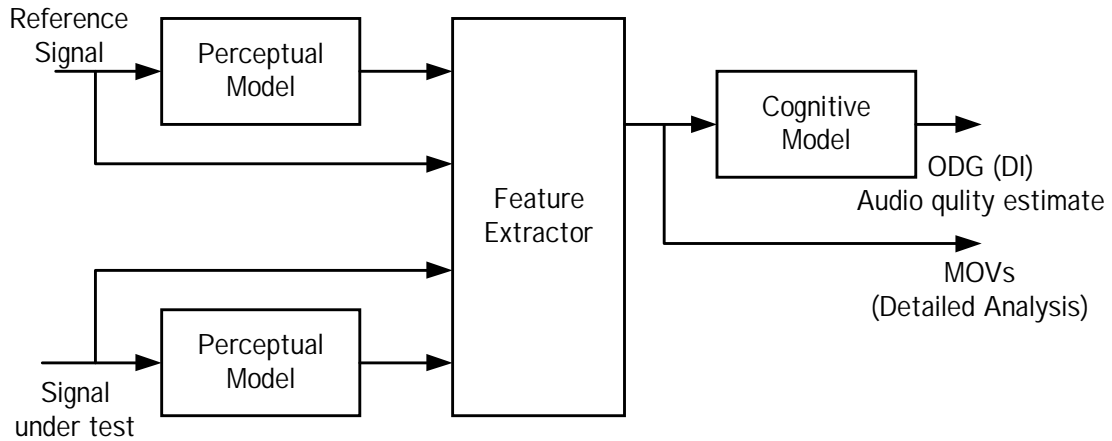


Figure B.2 – Representation of PEAQ model

## B.2 Basic version

The “Basic” version implements an FFT based ear model, as outlined in figure B.3.

Most features of this model are based on the fundamental psychoacoustic principles. Figure B.3 shows the signal flow from the input signal to the final calculation of the excitation pattern. The processing starts by a transformation of the input signal to the frequency domain. A 2048-point FFT is applied along with subsequent scaling of the spectra, according to the listening level, which has to be input by the user as a parameter. This results in the frequency resolution of approximately 23,4 Hz, and a corresponding temporal resolution of 23,4 ms (at 48 kHz sampling rate).

In the constructive block, the effects of the outer and middle ear are modelled by weighting the spectrum with the appropriate filter functions. Afterwards the spectra are grouped into critical bands, archiving a resolution of 1/4 bark per band. The subsequent adding of “internal noise” is intended to model effects, such as the permanent masking of sounds in our auditory system caused by the streaming of blood and other physiological phenomena. This step is followed by calculation of masking effects. Simultaneous masking is modelled by a frequency and level dependent spreading function. Temporal masking is modelled only partly since the temporal resolution is the same range as the timing of any background masking effects, which therefore cannot be modelled. Nevertheless, experiments have shown that backward masking is very coarsely modelled by side effects of the FFT.

Using the feature extractor, eleven MOVs are extracted from the compensation of the ear model output. Table B.1 shows a list of those MOVs and their interpretation. For further information about the MOVs please refer to the annex of the ITU-R recommendation BS.1387.

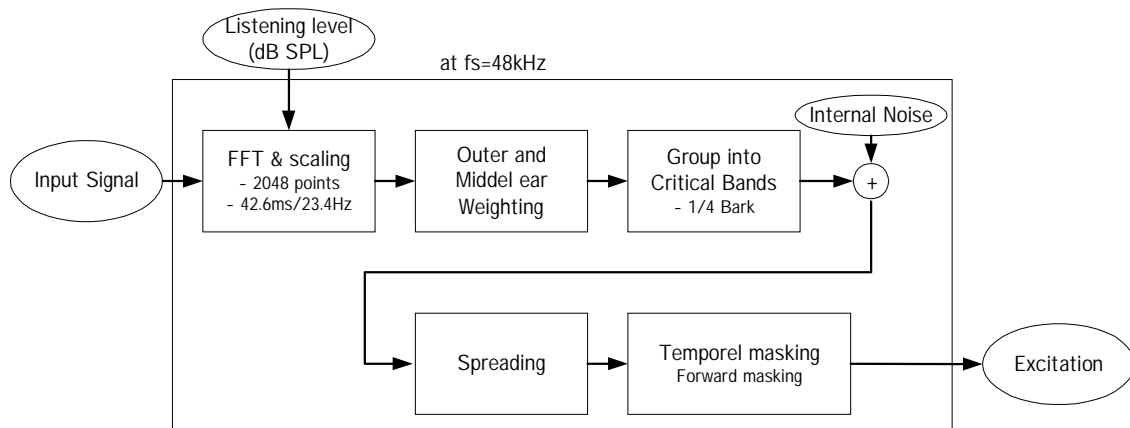


Figure B.3 – FFT based ear model, PEAQ basic version

**Table B.1 – Model output variables, PEAQ basic version**

Model Output Variable (MOV)	purpose
WinModDiff1 <sub>B</sub>	Changes in modulation (related to roughness)
AvgModDiff1 <sub>B</sub>	
AvgModDiff2 <sub>B</sub>	
RmsNoiseLoud <sub>B</sub>	Loudness of the distortion
BandwidthRef <sub>B</sub>	Linear distortions (frequency response etc.)
BandwidthTest <sub>B</sub>	
RelDistFrames <sub>B</sub>	Frequency of audible distortions
Total NMR <sub>B</sub>	Noise-to-mask ratio
MFPD <sub>B</sub>	Detection probability
ADB <sub>B</sub>	
EHS <sub>B</sub>	
	Harmonic structure of the error

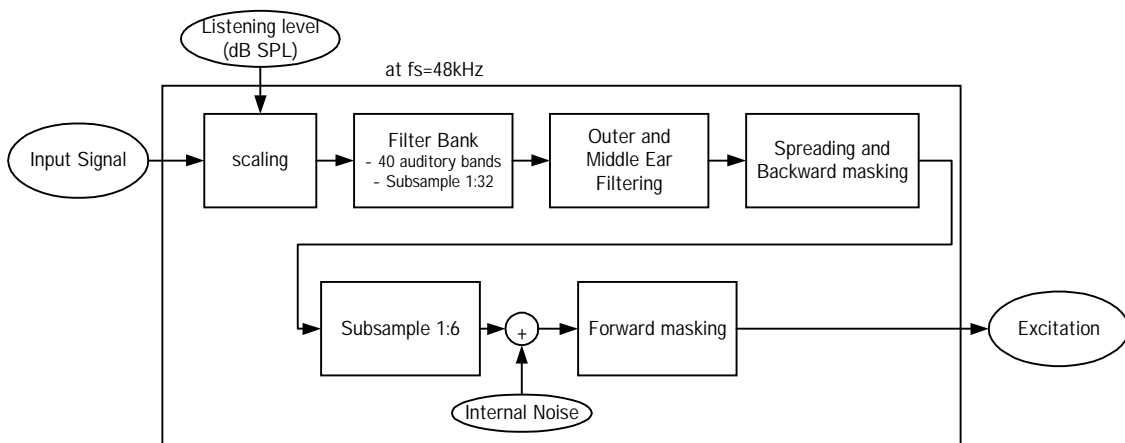
### B.3 Advanced version

The “Advanced” version use some MOVs derived by implementing the ear model of the “Basic” version but in addition to that it introduces a second ear model with improved temporal resolution, as illustrated in figure B.4.

Compared to the “Basic” version, this model performs the time frequency warping using a filter bank, thus grouping the signal into 40 auditory bands with a temporal resolution of approximately 0,66 ms. This allows for a very accurate modelling of backward masking effects. After the calculation of backward and simultaneous masking, the signal is sub-sampled by a factor of 1:6 in order to improve the computational efficiency. After adding the internal noise to the sub-sampled signal and finally modelling the forward masking effects, the output of this model is again the excitation.

In comparison to the FFT based “Basic” approach, the temporal resolution is improved, thus allowing for better simulation of temporal effects, at the cost of frequency resolution and computational complexity.

Due to the combination of parameters derived from both of the ear models, the number of MOVs used by the “Advanced” version to derive the final quality measure could be reduced to five, while simultaneously the accuracy of the algorithm was slightly improved compared to the “Basic” version. The MOVs used by the “Advanced” version are listed in table 5.2. For more detailed information about the advanced version, see the annex of the ITU-R recommendation BS.1387.

**Figure B.4 – Filter bank based ear model, PEAQ advanced version**

**Table B.2 – Model output variables, PEAQ advanced version**

<b>Model Output Variable (MOV)</b>	<b>Purpose</b>
RmsNoiseLoudAsym <sub>A</sub>	Loudness of the distortion
RmsModDiff <sub>A</sub>	Changes in modulation (related to roughness)
AvgLinDist <sub>A</sub>	Linear distortions (frequency response etc.)
Segmental NMR <sub>B</sub>	Noise-to-mask ratio
EHS <sub>B</sub>	Harmonic structure of the error

#### **B.4 Output value of PEAQ method**

The Objective Difference Grade (ODG) is the output value of PEAQ method that corresponds to the Subjective Difference Grade (SDG) in the subjective domain. The resolution of the ODG is limited to one decimal. However, one should be cautious and not generally expect that a difference between any pair of ODGs of tenth of a grade is significant. The same remark is valid when looking at results from a subjective listening test. The ODG can also assume positive values. Such values can occur because PEAQ use the cognitive model to map the MOVs to the results of subjective listening test. In the case of subjective listening tests, the SDG can assume a positive value, when a test person has incorrectly assigned the reference and test signal.

The Distortion Index (DI) has the same meaning as the ODG. However, DI and ODG can only be compared quantitatively but not qualitatively. The DI is characterized by a saturation that is less than the saturation of the ODG value. Furthermore, the range of values is different. As a general rule, you should use the ODG as the quality measure for ODG values greater than approximately  $-3,6$ . The ODG correlate very well with subjective assessment in this range. When ODG value is less than  $-3,6$  you should use the DI.

#### **B.5 Performance of PEAQ measurement method**

In order to validate the performance of PEAQ model, a number of different criteria may be relevant. The correlation between ODG and SDG is an obvious criterion to evaluate. In addition two further criteria that consider the reliability of the mean value were used for validation – the Absolute Error Score (AES) and the Tolerance Scheme.

The validation tests performed by ITU-R showed that PEAQ predicts the perceived quality with high-accuracy and is superior to previously existing measurements method. For further information please refer to the annex of the ITU-R recommendation BS.1387 and [AES-PEAQ]<sup>1</sup>.

<sup>1</sup> T. Theide et.al. "PEAQ – The ITU standard for Objective Measurement of Perceived Audio Quality," J. Audio Eng. Soc., vol.48, pp 3-29 (2000 Jan./Feb.)

## Bibliography

- ITU-T Recommendation P.930, Principles of a reference impairment system for video.
- ITU-T Recommendation G.113, Transmission impairments, Annex I: Provisional planning values for the equipment impairment factor.
- ITU-T Recommendation P.861, Objective quality measurement of telephone-band (300-3400 Hz) speech codecs.
- T. Theide et.al. “PEAQ – The ITU standard for Objective Measurement of Perceived Audio Quality,” J. Audio Eng. Soc., vol.48, pp 3-29 (2000 Jan./Feb.)
- Measuring quality in videoconferencing systems, Part number PC316, Intel Corporation (November 1997)
- Criteria for product evaluation, NASA Desktop video expert center, National Aeronautics and Space Administration, Ames Research Center, Moffett Field, California (August 1997)
- Quality aspects of computer-based video services, Norbert Gerfelder (Fraunhofer Institute for Computer Graphics, Darmstadt, Germany and Wolfgang Muller (Darmstadt Technical University), (Oct. 1995)
- Comparative study on narrow-bandwidth presentation of streaming educational videos, H. Ikeda, S. Dickerson, Y. Higaki, Journal of Faculty of Engineering, Chiba University, Vol.49, No. 1, pp.19-26 (1997-9).