



INTERNATIONAL TELECOMMUNICATION UNION

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**P.910**

(09/99)

SERIES P: TELEPHONE TRANSMISSION QUALITY,  
TELEPHONE INSTALLATIONS, LOCAL LINE  
NETWORKS

Audiovisual quality in multimedia services

---

**Subjective video quality assessment methods  
for multimedia applications**

ITU-T Recommendation P.910

(Previously CCITT Recommendation)

---

ITU-T P-SERIES RECOMMENDATIONS

**TELEPHONE TRANSMISSION QUALITY, TELEPHONE INSTALLATIONS, LOCAL LINE NETWORKS**

Vocabulary and effects of transmission parameters on customer opinion of transmission quality	Series P.10
Subscribers' lines and sets	Series P.30 P.300
Transmission standards	Series P.40
Objective measuring apparatus	Series P.50 P.500
Objective electro-acoustical measurements	Series P.60
Measurements related to speech loudness	Series P.70
Methods for objective and subjective assessment of quality	Series P.80 P.800
<b>Audiovisual quality in multimedia services</b>	<b>Series P.900</b>

*For further details, please refer to ITU-T List of Recommendations.*

**SUBJECTIVE VIDEO QUALITY ASSESSMENT METHODS  
FOR MULTIMEDIA APPLICATIONS**

**Summary**

This Recommendation describes non-interactive subjective assessment methods for evaluating the one-way overall video quality for multimedia applications such as videoconferencing, storage and retrieval applications, tele-medical applications, etc. These methods can be used for several different purposes including, but not limited to, selection of algorithms, ranking of audiovisual system performance and evaluation of the quality level during and audiovisual connection. This Recommendation also outlines the characteristics of the source sequences to be used, like duration, kind of content, number of sequences, etc.

**Source**

ITU-T Recommendation P.910 was revised by ITU-T Study Group 12 (1997-2000) and was approved under the WTSC Resolution No. 1 procedure on 30 September 1999.

## FOREWORD

ITU (International Telecommunication Union) is the United Nations Specialized Agency in the field of telecommunications. The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of the ITU. The ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Conference (WTSC), which meets every four years, establishes the topics for study by the ITU-T Study Groups which, in their turn, produce Recommendations on these topics.

The approval of Recommendations by the Members of the ITU-T is covered by the procedure laid down in WTSC Resolution No. 1.

In some areas of information technology which fall within ITU-T's purview, the necessary standards are prepared on a collaborative basis with ISO and IEC.

## NOTE

In this Recommendation, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

## INTELLECTUAL PROPERTY RIGHTS

The ITU draws attention to the possibility that the practice or implementation of this Recommendation may involve the use of a claimed Intellectual Property Right. The ITU takes no position concerning the evidence, validity or applicability of claimed Intellectual Property Rights, whether asserted by ITU members or others outside of the Recommendation development process.

As of the date of approval of this Recommendation, the ITU had not received notice of intellectual property, protected by patents, which may be required to implement this Recommendation. However, implementors are cautioned that this may not represent the latest information and are therefore strongly urged to consult the TSB patent database.

© ITU 2000

All rights reserved. No part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from the ITU.

## CONTENTS

	<b>Page</b>
1 Scope.....	1
2 References.....	1
3 Terms and definitions .....	2
4 Abbreviations.....	3
5 Source signal.....	3
5.1 Recording environment.....	3
5.2 Recording system.....	4
5.2.1 Camera.....	4
5.2.2 Video signal and storage format.....	4
5.3 Scene characteristics .....	4
5.3.1 Spatial perceptual information measurement .....	5
5.3.2 Temporal perceptual information measurement.....	5
6 Test methods and experimental design.....	5
6.1 Absolute Category Rating (ACR).....	6
6.2 Degradation Category Rating (DCR).....	7
6.3 Pair Comparison method (PC).....	7
6.4 Comparison of the methods .....	8
6.5 Reference conditions.....	9
6.6 Experimental design .....	9
7 Evaluation procedures.....	9
7.1 Viewing conditions .....	10
7.2 Processing and playback system .....	10
7.3 Viewers .....	10
7.4 Instructions to viewers and training session .....	11
8 Statistical analysis and reporting of results.....	11
Annex A – Details related to the characterization of the test sequences.....	12
A.1 Sobel filter.....	12
A.2 How to use SI and TI for test sequence selection .....	13
A.3 Examples.....	13
Annex B – Additional evaluative scales .....	15
B.1 Rating scales .....	14
B.2 Additional rating dimensions.....	16

	<b>Page</b>
Annex C – Simultaneous presentation of sequence pairs .....	18
C.1 Introduction.....	18
C.2 Synchronization .....	18
C.3 Viewing conditions .....	18
C.4 Presentations .....	19
Annex D – Video classes and their attributes .....	19
Appendix I – Bibliography.....	20
Appendix II – Test sequences .....	21
Appendix III – Instructions for viewing tests.....	22
III.1 ACR .....	22
III.2 DCR .....	22
III.3 PC.....	22
Appendix IV – The simultaneous double stimulus for a continuous evaluation .....	23
IV.1 Test procedure.....	23
IV.2 The training phase.....	23
IV.3 Test protocol features.....	23
IV.4 Data processing.....	24
IV.5 Reliability of the subjects.....	27
Appendix V – The object-based evaluation .....	28
Appendix VI – An additional evaluative scale for DRC.....	30

## Recommendation P.910

### SUBJECTIVE VIDEO QUALITY ASSESSMENT METHODS FOR MULTIMEDIA APPLICATIONS

(Geneva, 1996, 1999)

#### 1 Scope

This Recommendation is intended to define non-interactive subjective assessment methods for evaluating the quality of digital video images coded at bit rates specified in classes for TV3, MM4, MM5 and MM6, as specified in Table D.2 for applications such as videotelephony, videoconferencing and storage and retrieval applications. The methods can be used for several different purposes including, but not limited to, selection of algorithms, ranking of video system performance and evaluation of the quality level during a video connection.

#### 2 References

The following ITU-T Recommendations and other references contain provisions which, through reference in this text, constitute provisions of this Recommendation. At the time of publication, the editions indicated were valid. All Recommendations and other references are subject to revision; all users of this Recommendation are therefore encouraged to investigate the possibility of applying the most recent edition of the Recommendations and other references listed below. A list of the currently valid ITU-T Recommendations is regularly published.

- [1] ITU-T Recommendation P.930 (1996), *Principles of a reference impairment system for video*.
- [2] ITU-T Recommendation P.920 (1996), *Interactive test methods for audiovisual communications*.
- [3] ITU-R Recommendation BT.601-4 (1994), *Encoding parameters of digital television for studios*.
- [4] ITU-R Recommendation BT.500-9 (1998), *Methodology for the subjective assessment of the quality of television pictures*.
- [5] IEC Publication 60268-13, *Sound system equipment – Part 13: listening tests on loudspeakers*.
- [6] CCITT: *Handbook on Telephony*, Geneva, 1992.
- [7] ITU-T Recommendation P.800 (1996), *Methods for subjective determination of transmission quality*.
- [8] ITU-R Recommendation BT.814-1 (1993), *Specifications and alignment procedures for setting of brightness and contrast of displays*.
- [9] ITU-R Recommendation BT.815-1 (1993), *Specification of a signal for measurement of the contrast ratio of displays*.
- [10] CCIR Report 1213, *Test pictures and sequences for subjective assessments of digital codecs*, Annex to Volume XI, Part 1, 1990.
- [11] CCITT Recommendation J.61 (1990), *Transmission performance of television circuits designed for use in international connections*.

- [12] ITU-R Recommendation BT.812 (1992), *Subjective assessment of the quality of alphanumeric and graphic pictures in Teletext and similar services.*

### 3 Terms and definitions

This Recommendation defines the following terms:

- 3.1 gamma:** A parameter that describes the discrimination between the grey level steps on a visual display. The relation between the screen luminance and the input signal voltage is non-linear, with the voltage raised to an exponent gamma. To compensate for this non-linearity, a correction factor that is an inverse function of gamma is generally applied in the camera. Gamma also has an impact on colour rendition.
- 3.2 optimization tests:** Subjective tests that are typically carried out during either the development or the standardization of a new algorithm or system. The goal of these tests is to evaluate the performance of new tools in order to optimize the algorithms or the systems that are under study.
- 3.3 qualification tests:** Subjective tests that are typically carried out in order to compare the performance of commercial systems or equipment. These tests must be carried out under test conditions that are as much representative as possible of the real conditions of use.
- 3.4 spatial perceptual information (SI):** A measure that generally indicates the amount of spatial detail of a picture. It is usually higher for more spatially complex scenes. It is not meant to be a measure of entropy nor associated with the information defined in communication theory. See 5.3.1 for the equation for SI.
- 3.5 temporal perceptual information (TI):** A measure that generally indicates the amount of temporal changes of a video sequence. It is usually higher for high motion sequences. It is not meant to be a measure of entropy nor associated with the information defined in communication theory. See 5.3.2 for the equation for TI.
- 3.6 transparency (fidelity):** A concept describing the performance of a codec or a system in relation to an ideal transmission system without any degradation.
- Two types of transparency can be defined:
- The first type describes how well the processed signal conforms to the input signal, or ideal signal, using a mathematical criterion. If there is no difference, the system is fully transparent. The second type describes how well the processed signal conforms to the input signal, or ideal signal, for a human observer. If no difference can be perceived under any experimental condition, the system is perceptually transparent. The term "transparent" without explicit reference to a criterion will be used for systems that are perceptually transparent.
- 3.7 replication:** Repetition of the same circuit condition (with the same source material) for the same subject.
- 3.8 reliability of a subjective test:**
- intra-individual ("within subject") reliability refers to the agreement between a certain subject's repeated ratings of the same test condition;
  - inter-individual ("between subjects") reliability refers to the agreement between different subjects' ratings of the same test condition.
- 3.9 validity of a subjective test:** Agreement between the mean value of ratings obtained in a test and the true value which the test purports to measure.
- 3.10 reference conditions:** Dummy conditions added to the test conditions in order to anchor the evaluations coming from different experiments.
- 3.11 explicit reference (source reference):** The condition used by the assessors as reference to express their opinion, when the DCR method is used. This reference is displayed first within each



pair of sequences. Usually the format of the explicit reference is the format used at the input of the codecs under test (e.g.: ITU-R BT.601-4, CIF, QCIF, SIF, etc.). In the body of this Recommendation, the words "explicit" and "source" will be omitted whenever the context will make clear the meaning of "reference".

**3.12 implicit reference:** The condition used by the assessors as reference to express their opinion on the test material, when the ACR method is used. If the implicit reference is suggested by the experimenter, it must be well known to all the assessors (e.g. conventional TV systems, reality).

## 4 Abbreviations

This Recommendation uses the following abbreviations:

ACR	Absolute Category Rating
CCD	Charge Coupled Device
CI	Confidence Interval
CIF	Common Intermediate Format (picture format defined in Recommendation H.261 for video phone: 352 lines × 288 pixels)
CRT	Cathode Ray Tube
DCR	Degradation Category Rating
%GOB	Percent of Good or Better (proportion of Good and Excellent)
LCD	Liquid Crystal Display
MOS	Mean Opinion Score
PC	Pair Comparison
%POW	Percent of Poor or Worse (proportion of Poor and Bad votes)
QCIF	Quarter CIF (picture format defined in Recommendation H.261 for video phone: 176 lines × 144 pixels)
S/N	Signal-to-Noise ratio
SI	Spatial Information
SIF	Standard Intermediate Format [picture formats defined in ISO 11172 (MPEG-1): 352 lines × 288 pixels × 25 frames/s and 352 lines × 240 pixels × 30 frames/s]
SP	Simultaneous Presentation
std	Standard Deviation
TI	Temporal Information
VTR	Video Tape Recorder

## 5 Source signal

In order to control the characteristics of the source signal, the test sequences should be defined according to the goal of the test and recorded on a digital storage system. When the experimenter is interested in comparing results from different laboratories, it is necessary to use a common set of source sequences to eliminate a further source of variation.

### 5.1 Recording environment

Lighting source(s) (bulbs or fluorescent lamps) can be placed above or on the side of the camera. When placing the lights, recognize that overhead is more typical of office lighting, and should be used with scenes that portray the business environment. Studio lights and other non-typical sources should be avoided.

The lighting conditions of the room in the field of view could vary from 100 lux to about 10 000 lux for indoor use. The variation (AC frequency) of the light (fluorescent lighting) must be taken into account because this may cause a flicker in the recorded video sequence.

Lighting conditions, wall colours, surface reflectance, etc. should be carefully controlled and reported.

## **5.2 Recording system**

### **5.2.1 Camera**

Picture sequences should be recorded by a high quality CCD camera.

The signal-to-noise ratio of the input video signal can strongly affect the performance of the codec.

To define the video input, the following points should be specified:

- the dynamic range of the Y U V signals;
- the gamma correction factor (should be 0.45);
- the bandwidth/slopes of the filters;
- the sensitivity of the camera at very low lighting conditions and the characteristics of an Automatic Gain Control (AGC), if used.

The weighted S/N should be measured according to Recommendation J.61 Part C, subclause 3.2.1 [11]. The weighted S/N should be greater than 45 dB r.m.s.

The instability or the jitters of the clock signals could cause noise effects. A minimum stability of 0.5 ppm is required for the camera clocking device.

Either fixed or variable focal length systems can be used. For desktop terminals a focal depth from 30 cm to 120 cm is reasonable, while for multi-user systems a focal depth from 50 cm to infinity might be more appropriate. To support the variation of illuminance in the recording room, either an adjustable iris or neutral density filters should be used. The camera should have an automatic white balance so that adaptation to the colour temperature of the light source can be accomplished. The correction of white temperature can range from 2700° K (indoor use with electrical bulb) to 6500° K (daylight temperature with clouded sky).

### **5.2.2 Video signal and storage format**

Video source signals provided by the camera should be sampled in conformance with Part A of [3]. In order to avoid distortion of the source signal, it should be stored in digital format, e.g. on computer or D1 4:2:2 tape format.

## **5.3 Scene characteristics**

The selection of test scenes is an important issue. In particular, the spatial and temporal perceptual information of the scenes are critical parameters. These parameters play a crucial role in determining the amount of video compression that is possible, and consequently, the level of impairment that is suffered when the scene is transmitted over a fixed-rate digital transmission service channel. Fair and relevant video test scenes must be chosen such that their spatial and temporal information is consistent with the video services that the digital transmission service channel was intended to provide. The set of test scenes should span the full range of spatial and temporal information of interest to users of the devices under test.

Details on the characterization of the test sequences and examples of suitable test scenes are given in Annex A and in Appendices II and III.

The number of sequences should be defined according to the experimental design. In order to avoid boring the observers and to achieve a minimum reliability of the results, at least four different types of scenes (i.e. different subject matters) should be chosen for the sequences.

The following subclauses present methods for quantifying the spatial and temporal information of test scenes. These methods for evaluating the spatial and temporal information of test scenes are applicable to video quality testing both now and in the future. The location of the video scene within the spatial-temporal matrix is important because the quality of a transmitted video scene (especially after passing through a low bit-rate codec) is often highly dependent on this location. The spatial and temporal information measures presented here can be used to assure appropriate coverage of the spatial-temporal plane.

The spatial and temporal information measures given below are single-valued for each frame over a complete test sequence. This results in a time series of values which will generally vary to some degree. The perceptual information measures given below remove this variability with a maximum function (maximum value for the sequence). The variability itself may be usefully studied for example with plots of spatial-temporal information on a frame-by-frame basis. The use of information distributions over a test sequence also permits better assessment of scenes with scene cuts.

### 5.3.1 Spatial perceptual information measurement

The Spatial perceptual Information, SI, is based on the Sobel filter. Each video frame (luminance plane) at time  $n$  ( $F_n$ ) is first filtered with the Sobel filter [ $Sobel(F_n)$ ]. The standard deviation over the pixels ( $std_{space}$ ) in each Sobel-filtered frame is then computed. This operation is repeated for each frame in the video sequence and results in a time series of spatial information of the scene. The maximum value in the time series ( $max_{time}$ ) is chosen to represent the spatial information content of the scene. This process can be represented in equation form as:

$$SI = \max_{time} \{std_{space}[Sobel(F_n)]\}$$

### 5.3.2 Temporal perceptual information measurement

The Temporal perceptual Information, TI, is based upon the motion difference feature,  $M_n(i, j)$ , which is the difference between the pixel values (of the luminance plane) at the same location in space but at successive times or frames.  $M_n(i, j)$  as a function of time ( $n$ ) is defined as:

$$M_n(i, j) = F_n(i, j) - F_{n-1}(i, j)$$

here  $F_n(i, j)$  is the pixel at the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column of  $n^{\text{th}}$  frame in time.

The measure of Temporal Information, TI, is computed as the maximum over time ( $max_{time}$ ) of the standard deviation over space ( $std_{space}$ ) of  $M_n(i, j)$  over all  $i$  and  $j$ .

$$TI = \max_{time} \{std_{space}[M_n(i, j)]\}$$

More motion in adjacent frames will result in higher values of TI.

NOTE – For scenes that contain scene cuts, two values may be given: one where the scene cut is included in the temporal information measure, and one where it is excluded from the measurement.

## 6 Test methods and experimental design

Measurement of the perceived quality of images requires the use of subjective scaling methods. The condition for such measurements to be meaningful is that there exists a relation between the physical characteristics of the "stimulus", in this case the video sequence presented to the subjects in a test, and the magnitude and nature of the sensation caused by the stimulus.

A number of experimental methods have been validated for different purposes. Here three methods are recommended for applications using connections at bit rates specified in classes for TV3, MM4, MM5 and MM6, as specified in Table D.2. Further test methods are described in Appendices IV and V.

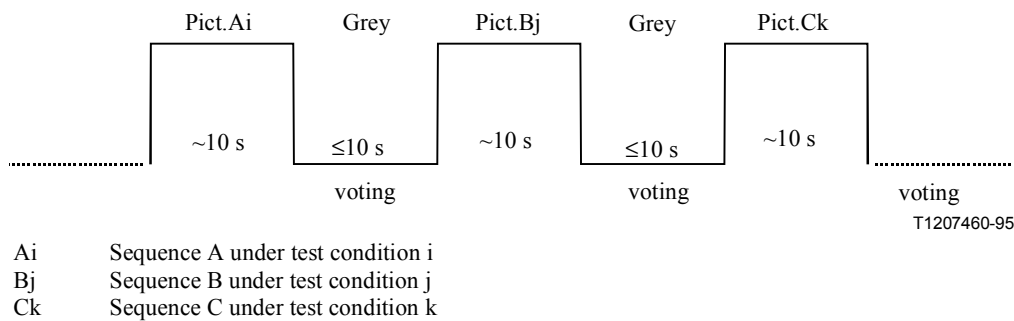
The final choice of one of these methods for a particular application depends on several factors, such as the context, the purpose and where in the development process the test is to be performed.

### 6.1 Absolute Category Rating (ACR)

The Absolute Category Rating method is a category judgement where the test sequences are presented one at a time and are rated independently on a category scale. (This method is also called Single Stimulus Method.)

The method specifies that after each presentation the subjects are asked to evaluate the quality of the sequence shown.

The time pattern for the stimulus presentation can be illustrated by Figure 1. The voting time should be less than or equal to 10 s, depending upon the voting mechanism used. The presentation time may be reduced or increased according to the content of the test material.



**Figure 1/P.910 – Stimulus presentation in the ACR method**

The following five-level scale for rating overall quality should be used:

- 5      Excellent
- 4      Good
- 3      Fair
- 2      Poor
- 1      Bad

If higher discriminative power is required, a nine-level scale may be used. Examples of suitable numerical or continuous scales are given in Annex B. Annex B also gives examples of rating dimensions other than overall quality. Such dimensions may be useful for obtaining more information on different perceptual quality factors when the overall quality rating is nearly equal for certain systems under test, although the systems are clearly perceived as different.

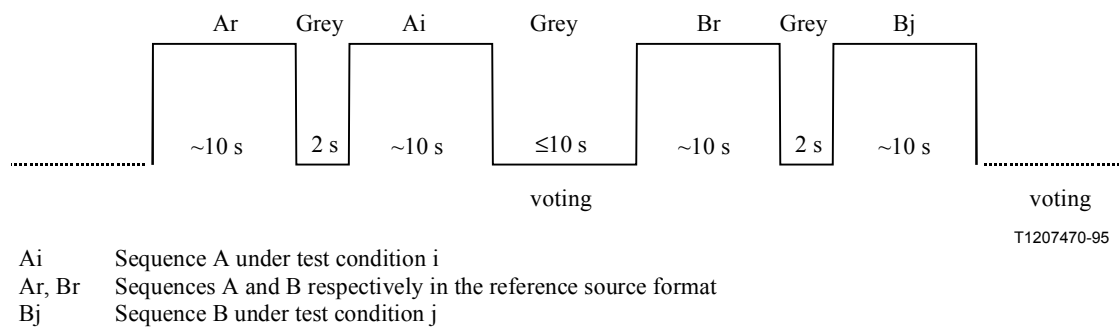
For the ACR method, the necessary number of replications is obtained by repeating the same test conditions at different points of time in the test.

## 6.2 Degradation Category Rating (DCR)

The Degradation Category Rating implies that the test sequences are presented in pairs: the first stimulus presented in each pair is always the source reference, while the second stimulus is the same source presented through one of the systems under test. (This method is also called the Double Stimulus Impairment Scale method.)

When reduced picture formats are used (e.g. CIF, QCIF, SIF), it could be useful to display the reference and the test sequence simultaneously on the same monitor. Guidelines on this presentation procedure are discussed in Annex C.

The time pattern for the stimulus presentation can be illustrated by Figure 2. The voting time should be less than or equal to 10 s, depending upon the voting mechanism used. The presentation time may be reduced or increased according to the content of the test material.



**Figure 2/P.910 – Stimulus presentation in the DCR method**

In this case the subjects are asked to rate the impairment of the second stimulus in relation to the reference.

The following five-level scale for rating the impairment should be used:

- 5    Imperceptible
- 4    Perceptible but not annoying
- 3    Slightly annoying
- 2    Annoying
- 1    Very annoying

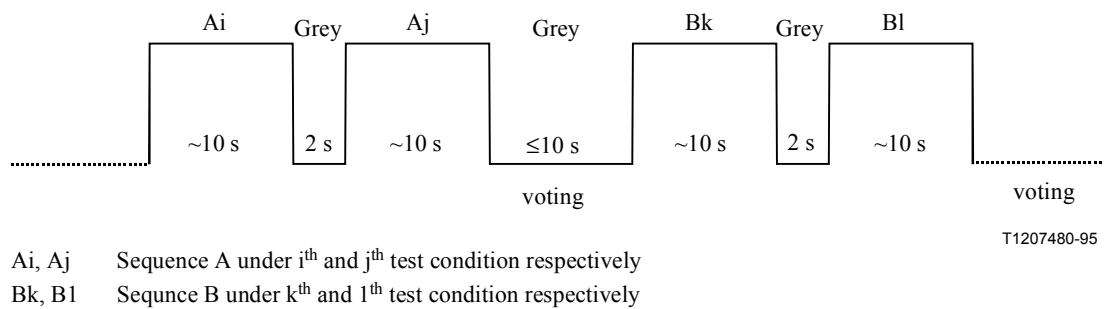
The necessary number of replications is obtained for the DCR method by repeating the same test conditions at different points of time in the test.

## 6.3 Pair Comparison method (PC)

The method of Pair Comparisons implies that the test sequences are presented in pairs, consisting of the same sequence being presented first through one system under test and then through another system.

The systems under tests (A, B, C, etc.) are generally combined in all the possible  $n(n-1)$  combinations AB, BA, CA, etc. Thus, all the pairs of sequences should be displayed in both the possible orders (e.g. AB, BA). After each pair a judgement is made on which element in a pair is preferred in the context of the test scenario.

The time pattern for the stimulus presentation can be illustrated by Figure 3. The voting time should be less than or equal to 10 s, depending upon the voting mechanism used. The presentation time should be about 10 s and it may be reduced or increased according to the content of the test material.



**Figure 3/P.910 – Stimulus presentation in the PC method**

When reduced resolutions are used (e.g. CIF, QCIF, SIF), it could be useful to display each pair of sequences simultaneously on the same monitor. Guidelines on this presentation procedure are discussed in Annex C.

For the PC method, the number of replications need not generally be considered, because the method itself implies repeated presentation of the same conditions, although in different pairs.

A variation of the PC method utilizes a categorical scale to further measure the differences between the pair of sequences. See References [4] and [7].

#### 6.4 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) and methods that do not use any explicit reference (e.g. ACR and PC). This second class of method does not test transparency or fidelity.

The DCR method should be used when testing the fidelity of transmission with respect to the source signal. This is frequently an important factor in the evaluation of high quality systems. DCR has long been a key method specified in [4], for the assessment of television pictures whose typical quality represents the extreme high levels of videotelephony and videoconferencing. Other methods may also be used to evaluate high quality systems. The specific comments of the DCR scale (imperceptible/perceptible) are valuable when the viewer's detection of impairment is an important factor.

Thus, when it is important to check the fidelity with respect to the source signal, DCR method should be used.

DCR should also be applied for high quality system evaluation in the context of multimedia communication. Discrimination of imperceptible/perceptible impairment in the DCR scale supports this, as well as comparison with the reference quality.

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.

The principal merit of the PC method is its high discriminatory power, which is of particular value when several of the test items are nearly equal in quality.

When a large number of items are to be evaluated in the same test, the procedure based on the PC method tends to be lengthy. In such a case an ACR or DCR test may be carried out first with a limited number of observers, followed by a PC test solely on those items which have received about the same rating.

## 6.5 Reference conditions

The results of quality assessments often depend not only on the actual video quality, but also on other factors such as the total quality range of the test conditions, the experience and expectations of the assessors, etc. In order to control some of these effects, a number of dummy test conditions can be added and used as references.

A description of reference conditions and procedures to produce them is given in Recommendation P.930 [1]. The introduction of the source signal as a reference condition in a PC test is specially recommended when the impairments introduced by the test items are small.

The quality level of the reference conditions should cover at least the quality range of the test items.

## 6.6 Experimental design

Different experimental designs, such as complete randomized design, Latin, Graeco-Latin and Youden square designs, replicated block designs, etc. [I.5] can be used, the selection of which should be driven by the purpose of the experiment.

It is left to the experimenter to select a design method in order to meet specific cost and accuracy objectives. The design may also depend upon which conditions are of particular interest in a given test.

It is recommended to include at least two, if possible three or four, replications (i.e. repetitions of identical conditions) in the experiment. There are several reasons for using replications, the most important being that "within subject variation" can be measured using the replicated data. For testing the reliability of a subject the same order of presentation under identical conditions can be used. If a different order of presentation is used, the resulting variation in the experimental data is composed of the order effect and the within subject variation.

Replications make it possible to calculate individual reliability per subject and, if necessary, to discard unreliable results from some subjects. An estimate of both within- and between- subject standard deviation is furthermore a prerequisite for making a correct analysis of variance and to generalize results to a wider population. In addition, learning effects within a test are to some extent balanced out.

A further improvement in the handling of learning effects is obtained by including a training session in which at least five conditions are presented at the beginning of each test session. These conditions should be chosen to be representative of the presentations to be shown later during the session. The preliminary presentations are not to be taken into account in the statistical analysis of the test results.

## 7 Evaluation procedures

Table 1 lists typical viewing conditions as used in video quality assessment. The actual parameter settings used in the assessment should be specified. For the comparison of test results, all viewing conditions must be fixed and equal over laboratories for the same kind of tests.

Both the size and the type of monitor used should be appropriate for the application under investigation. When sequences are presented through a PC-based system, the characteristics of the display must be specified, e.g. dot pitch of the monitor, type of video display card used, etc.

Concerning the displaying format, it is preferable to use the whole screen for displaying the sequences. Nevertheless when, for some reason, the sequences must be displayed on a window of the screen, the colour of the background in the screen should be 50% grey corresponding to  $Y=U=V=128$  (U and V unsigned).

## 7.1 Viewing conditions

The test should be carried out under the following viewing conditions:

**Table 1/P.910 – Viewing conditions**

Parameter	Setting
Viewing distance (Note 1)	1-8 H (Note 2)
Peak luminance of the screen	100-200 cd/m (Note2)
Ratio of luminance of inactive screen to peak luminance	$\leq 0.05$
Ratio of the luminance of the screen, when displaying only black level in a completely dark room, to that corresponding to peak white	$\leq 0.1$
Ratio of luminance of background behind picture monitor to peak luminance of picture (Note 3)	$\leq 0.2$
Chromaticity of background (Note 4)	D <sub>65</sub>
Background room illumination (Note 3)	$\leq 20$ lux
NOTE 1 – For a given screen height, it is likely that the viewing distance preferred by the subjects increases when visual quality is degraded. Concerning this point, the preferred viewing distance should be predetermined for qualification tests. Viewing distance in general depends on the applications.	
NOTE 2 – H indicates the picture height. The viewing distance should be defined taking into account not only the screen size, but also the type of screen, the type of application and the goal of the experiment.	
NOTE 3 – This value indicates a setting allowing maximum detectability of distortions, for some applications higher values are allowed or they are determined by the application.	
NOTE 4 – For PC monitors the chromaticity of background may be adapted to the chromaticity of the monitor.	

## 7.2 Processing and playback system

There are two methods for obtaining test images from the source recordings:

- a) by transmitting or replaying the video recordings in real time through the systems under test, while subjects are watching and responding;
- b) by off-line processing of the source recordings through the device under test and recording the output to give a new set of recordings.

In the second case a digital VTR should be used to minimize the impairments that can be produced by the recording process. In any case, taking into account that the impairments introduced by low bit-rate coding schemes are usually more evident than the impairments introduced by modulation, professional quality VTRs such as D2, MII and BetacamSP can be used.

Either a CRT or an LCD monitor may be used. Both the size and the type of monitor used should be appropriate for the application under investigation.

The monitors should be aligned according to the procedures defined in [8].

## 7.3 Viewers

The possible number of subjects in a viewing test (as well as in usability tests on terminals or services) is from 4 to 40. Four is the absolute minimum for statistical reasons, while there is rarely any point in going beyond 40.



The actual number in a specific test should really depend on the required validity and the need to generalize from a sample to a larger population.

In general, at least 15 observers should participate in the experiment. They should not be directly involved in picture quality evaluation as part of their work and should not be experienced assessors.

Nevertheless, in the early phases in the development of video communication systems and in pilot experiments carried out before a larger test, small groups of experts (4-8) or other critical subjects can provide indicative results.

Prior to a session, the observers should usually be screened for normal visual acuity or corrected-to-normal acuity and for normal colour vision. Concerning acuity, no errors on the 20/30 line of a standard eye chart [I.3] should be made. The chart should be scaled for the test viewing distance and the acuity test performed at the same location where the video images will be viewed (i.e. lean the eye chart up against the monitor) and have the subjects seated. Concerning colour, no more than 2 plates [I.4] should be missed out of 12.

#### **7.4 Instructions to viewers and training session**

Before starting the experiment, a scenario of the intended application of the system under test should be given to the subjects. In addition, a description of the type of assessment, the opinion scale and the presentation of the stimuli is given in written form. The range and type of impairments should be presented in preliminary trials, which may contain video sequences other than those used in the actual tests.

It must not be implied that the worst quality seen in the training set necessarily corresponds to the lowest subjective grade on the scale.

Questions about procedure or about the meaning of the instructions should be answered with care to avoid bias and only before the start of the session.

A possible text for instructions to be given to the assessors is suggested in Appendix III.

### **8 Statistical analysis and reporting of results**

The results should be reported along with the details of the experimental set-up. For each combination of the test variables, the mean value and the standard deviation of the statistical distribution of the assessment grades should be given.

From the data, subject reliability should be calculated and the method used to assess subject reliability should be reported. Some criteria for subjective reliability are given in [4] and [5].

It is informative to analyse the cumulative distribution of scores. Since the cumulative distributions are not sensitive to linearity, these may be particularly useful for data for which the linearity is doubtful, as those obtained by using the ACR and DCR methods, together with category scales without grading (i.e. category judgement).

The data can be organized for example as shown in Table 2 for ACR.

**Table 2/P.910 – Informative table with cumulative distribution of scores for ACR method**

Condition	Total votes	Excellent	Good	Fair	Poor	Bad	MOS	CI	Std	%GOB	%POW

**Condition:** Label indicating a combination of test variables.  
**Total votes:** Number of votes collected for that condition.  
**Excellent, Fair ... Bad:** Occurrence of each vote.

The classical techniques of analysis of variance should be used to evaluate the significance of the test parameters. If the assessment is aimed at evaluating the video quality as a function of a parameter, curve fitting techniques can be useful for the interpretation of the data.

In the case of pair comparisons, the method to calculate the position of each stimulus on an interval scale, where the difference between the stimuli corresponds to the difference in preference, is described in the *Handbook on Telephony*, Section 2.6.2C [6].

## ANNEX A

### Details related to the characterization of the test sequences

#### A.1 Sobel filter

The Sobel filter is implemented by convolving two  $3 \times 3$  kernels over the video frame and taking the square root of the sum of the squares of the results of these convolutions.

For  $y = \text{Sobel}(x)$ , let  $x(i, j)$  denote the pixel of the input image at the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column.  $Gv(i, j)$  will be the result of the first convolution and is given as:

$$Gv(i, j) = -1 \times x(i-1, j-1) - 2 \times x(i-1, j) - 1 \times x(i-1, j+1) + \\ + 0 \times x(i, j-1) + 0 \times x(i, j) + 0 \times x(i, j+1) + \\ + 1 \times x(i+1, j-1) + 2 \times x(i+1, j) + 1 \times x(i+1, j+1)$$

Similarly,  $Gh(i, j)$  will be the result of the second convolution and is given as:

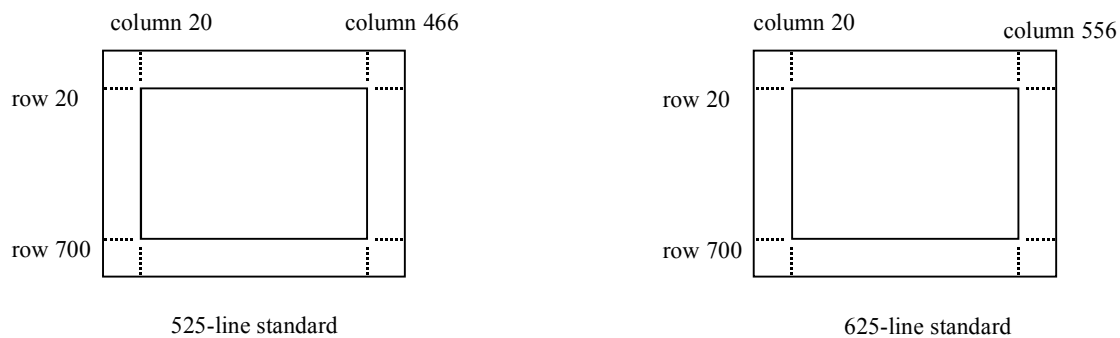
$$Gh(i, j) = -1 \times x(i-1, j-1) + 0 \times x(i-1, j) + 1 \times x(i-1, j+1) + \\ - 2 \times x(i, j-1) + 0 \times x(i, j) + 2 \times x(i, j+1) + \\ - 1 \times x(i+1, j-1) + 0 \times x(i+1, j) + 1 \times x(i+1, j+1)$$

Hence, the output of the Sobel filtered image at the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column is given as:

$$y(i, j) = \sqrt{[Gv(i, j)]^2 + [Gh(i, j)]^2}$$

The calculations are performed for all  $2 \leq i \leq N - 1$  and  $2 \leq j \leq M - 1$ , where  $N$  is the number of rows and  $M$  is the number of columns.

It is recommended that the calculations be performed on a subimage of the video frame to avoid unwanted edge effects and because the extreme edges of a video frame are usually invisible to CRT users. This can be accomplished by using a suitable subimage as illustrated for example in Figure A.1 for the 625- and 525-lines ITU-R BT.601-4 formats [3].



T1207490-95

**Figure A.1/P.910 – Subimages to be used to calculate SI and TI for 525- and 625-line ITU-R BT.601-4 formats [3]**

Further information on the Sobel filter can be found in [I.1].

## A.2 How to use SI and TI for test sequence selection

When selecting test sequences, it can be useful to compare the relative spatial information and temporal information found in the various sequences available. Generally, the compression difficulty is directly related to the spatial and temporal information of a sequence.

If a small number of test sequences are to be used in a given test, it may be important to choose sequences that span a large portion of the spatial-temporal information plane (see Figure A.2). In the case where four test sequences are to be used in a test, one might wish to choose a sequence from each of the four quadrants of the spatial-temporal information plane.

Alternately, if one were trying to choose test sequences which were equivalent in coding difficulty, then choosing sequences that had similar SI and TI values would be desirable.

## A.3 Examples

Figure A.2 shows the relative amounts of spatial and temporal information for some representative test scenes and how they can be placed on a spatial-temporal information plane.

Along the  $TI=0$  axis (along the bottom of the plot) are found the still scenes and those with very limited motion (such as l, f, and a). Near the top of the plot are found scenes with a lot of motion (such as p, q and i). Along the  $SI=0$  axis (at the left edge of the plot) are found scenes with minimal spatial detail (such as l, k, x, u and f). Near the right edge of the plot are found scenes with the most spatial detail (such as h and s). The values of SI and TI were obtained using the above equations and video which has been spatially sampled according to ITU-R BT.601-4 specifications [3]. Table A.1 lists the example test scenes by scene content category.

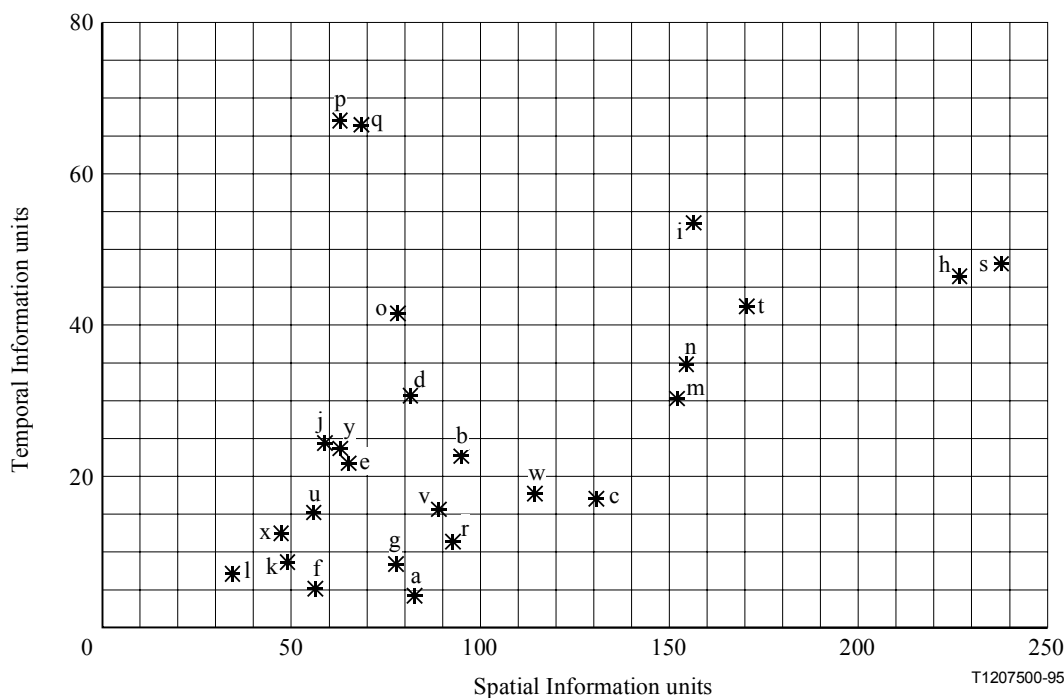


Figure A.2/P.910 – Spatial-temporal plot for example test scene set

Table A.1/P.910 – Scene content categories

Category	Description	Scene name and letter
A	One person, mainly head and shoulders, limited detail and motion	vtc1nw(f), susie(j), disguy(k), disgal(l)
B	One person with graphics and/or more detail	vtc2mp(a), vtc2zm(b), boblec(e), smity1(m), smity2(n), vowels(w), inspec(x)
C	More than one person	3inrow(d), 5row1(g), intros(o), 3twos(p), 2wbord(q), split6(r)
D	Graphics with pointing	washdc(c), cirkit(s), rodmap(t), filter(u), ysmite(v),
E	High object and/or camera motion (examples of broadcast TV)	flogar(h), ftball(i), fedas(y)

## ANNEX B

### Additional evaluative scales

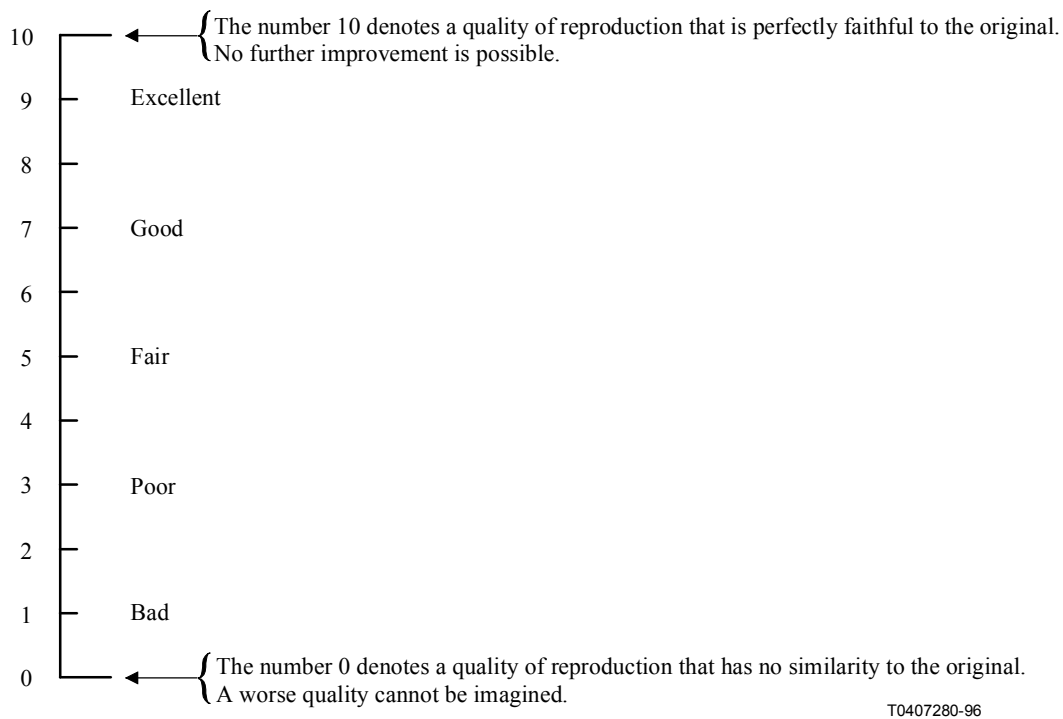
#### B.1 Rating scales

Particularly for the assessment of low bit-rate video codecs it is often necessary to use rating scales with more than five grades. A suitable scale for this purpose is the nine-grade scale, where the five verbally defined quality categories as recommended in 6.1 are used as labels for every second grade on the scale, as shown in Figure B.1.

9	Excellent
8	
7	Good
6	
5	Fair
4	
3	Poor
2	
1	Bad

**Figure B.1/P.910 – Nine-grade numerical quality scale**

A further extension of this scale is shown in Figure B.2, where the endpoints have been verbally defined as anchoring points which are not used for the rating. In this verbal definition some kind of reference is used (for example in Figure B.2 the original is used as reference). This reference can be either explicit or implicit and it will be clearly illustrated during the training phase. See also [5] and [6] Section 2.6 Scale a).



T0407280-96

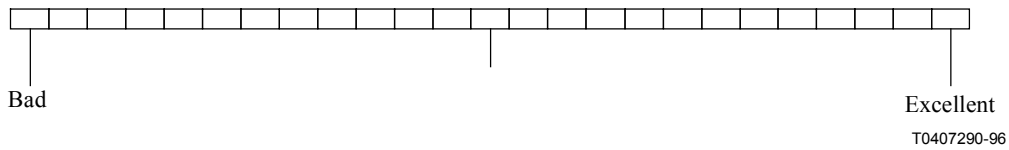
**Figure B.2/P.910 – Eleven-grade numerical quality scale**

For both types of scales the response from the subjects may be recorded either as numbers, which are written down on a response sheet, or as marks on the scale itself (in which case a separate scale has to be given on the response sheet for each rating condition). When numerical responses are required, the subjects should be encouraged to use decimals (e.g. 2.2 instead of 2) but they may still have the choice only to use integers.

It should be noted that it may be difficult to translate the names of the scale categories into different languages. In doing so the inter-category relationship could become different from that in the original language [I.6].

An additional possibility is to use continuous scales.

Since continuous data is usually rounded to some reasonable precision, to simplify data collection a voting scale like that shown in Figure B.3, can be used. Labels are used only at the endpoints and a mark is indicated in the middle of the scale. This should reduce the bias due to the interpretation of the labels. Each area can correspond to a specific numerical value and the data can be collected without ambiguity.



**Figure B.3/P.910 – Quasi-continuous scale for quality ratings**

## **B.2 Additional rating dimensions**

If the systems which are assessed in a test are judged as being rather equal in overall quality and therefore get very similar scores, it may be advantageous to rate additional quality components on separate scales for each condition. In this way it is possible to receive information on specific characteristics where the test objects are perceived as significantly different, even if the overall quality is in fact almost the same. Results from such additional tests can give valuable diagnostic information on the systems under test.

Examples of rating dimensions which may be assumed to define factors that contribute to the perceived global image quality are listed below, together with an indication of whether a factor contributes positively or negatively to quality:

- Brightness (positive);
- Contrast (positive);
- Colour reproduction (positive);
- Outline definition (positive);
- Background stability (positive);
- Speed in image reassembling (positive);
- Jerkiness (negative);
- "Smearing" effects (negative);
- "Mosquito" effects (negative);
- Double images/shadows (negative);
- Halo (negative).

Recent research has shown that these factors may be combined into a predicted global quality by giving appropriate weightings to each factor and then adding them together [I.2].

To evaluate separately the dimensions of the overall video quality, a special questionnaire can be used. Examples of questions that may be asked after the presentation of each test condition are given in the questionnaire below.

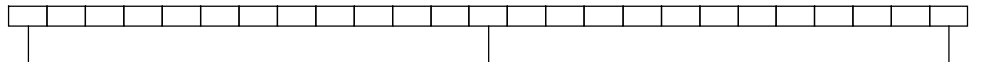
## Questionnaire

Could you kindly answer the following questions about the last sequence shown?

You can express your opinion by inserting a mark on the scales below.

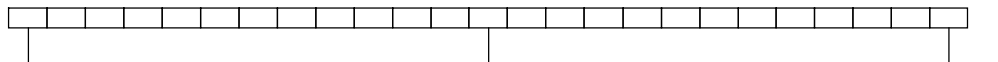
1) How would you rate image colours?

Bad Excellent  
T0407290-96



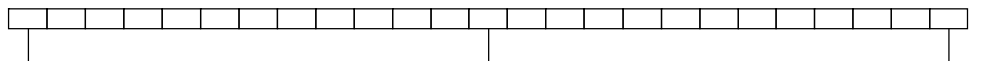
2) How would you rate image contrast?

Bad Excellent  
T0407290-96



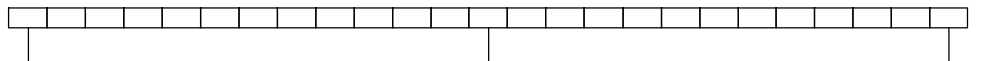
3) How would you rate the image borders?

Bad Excellent  
T0407290-96



4) How would you rate the movement continuity?

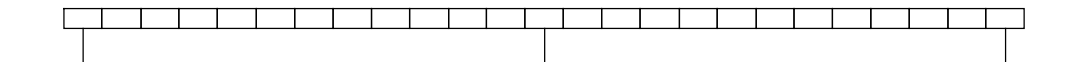
Bad Excellent  
T0407290-96



5) Did you notice any flicker in the sequence?  Yes  No

If you noticed flicker, please rate it on the scale below

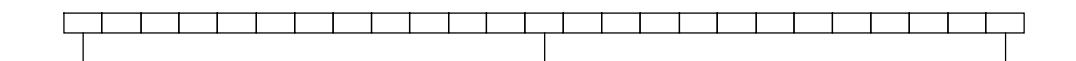
Very Annoying Not Annoying  
T0414170-00



6) Did you notice any smearing in the sequence?  Yes  No

If you noticed smearing, please rate it on the scale below

Very Annoying Not Annoying  
T0414170-00



NOTE – When these scales are used, all the quality/impairment categories taken into account (e.g. movement continuity, flicker, smearing, etc.) must be carefully illustrated during the training sessions.

## ANNEX C

### Simultaneous presentation of sequence pairs

#### C.1 Introduction

When the systems which are assessed in a test use reduced picture format, like CIF, QCIF, SIF, etc., and either the DCR or the PC methods are used, it may be advantageous to display simultaneously the two sequences of each pair on the same monitor.

The advantages in using Simultaneous Presentation (SP) are:

- 1) SP reduces considerably the duration of the test.
- 2) If suitable picture dimensions are used, it is easier for the subjects to evaluate the differences between the stimuli.
- 3) Since under the same test conditions the number of presentations is halved, the attention of the subjects is usually higher when the SP is used.

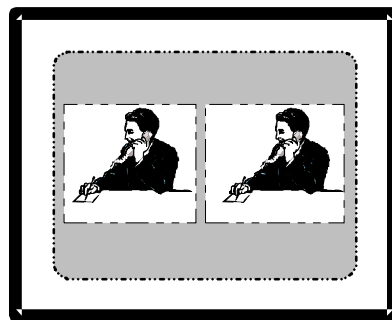
SP requires particular precautions in order to allow the subjects to avoid bias due to the type of presentation.

#### C.2 Synchronization

The two sequences must be perfectly synchronized; that means that they both must start and stop at the same frame and that the displaying must be synchronized. This does not preclude that sequences coded at different bit rates may be compared, provided that a suitable temporal up-sampling is applied.

#### C.3 Viewing conditions

The sequences must be displayed in two windows put side-by-side within a 50% grey background (the grey is specified in 5.1), as shown in Figure C.1. In order to reduce the eye movement to switch the attention between the two windows, the viewing distance should be  $8H$ , where  $H$  indicates the picture height. The diagonal dimension of the monitors should be at least 14 inches.



T1207510-95

Figure C.1/P.910 – Relative position of the two sequences in SP



#### C.4 Presentations

In DCR the reference should be placed always on the same side (e.g. left) and the subjects must be aware of the relative positions of reference and test conditions.

In PC all the pairs of sequences must be displayed in both the possible orders (e.g. AB, BA). This means that the sequences that were displayed on the left side are now displayed on the right one and vice versa.

### ANNEX D

#### Video classes and their attributes

In this Recommendation the highest video quality considered is ITU-R Recommendation BT.601, 8 bit/pixel linear PCM coded video in 4:2:2, Y, C<sub>R</sub>, C<sub>B</sub> format.

**Table D.1/P.910 – Definitions of video classes**

TV 0	Loss-less: ITU-R Rec. BT.601, 8 bit per pixel, video used for applications without compression.
TV 1	Used for complete post production, many edits and processing layers, intra-plant transmission. Also used for remote site to plant transmission. Perceptually transparent when compared to TV 0.
TV 2	Used for simple modifications, few edits, character/logo overlays, program insertion, and inter-facility transmission. A broadcast example would be network to affiliate transmission. Other examples are a cable system regional downlink to a local head-end and a high quality videoconferencing system. Nearly perceptually transparent when compared to TV 0.
TV 3	Used for delivery to home/consumer (no changes). Other examples are a cable system from the local head-end to a home and medium to high quality videoconferencing. Low artefacts are present when compared to TV 2.
MM 4	All frames encoded. Low artefacts relative to TV 3. Medium quality videoconferencing. Usually $\geq 30$ fps.
MM 5	Frames may be dropped at encoder. Perceivable artefacts possible, but quality level useful for designed tasks, e.g. low quality videoconferencing.
MM 6	Series of stills. Not Intended to provide full motion (Examples: Surveillance, Graphics).

**Table D.2/P.910 – Attributes of video classes**

<b>Video class</b>	<b>Spatial format</b>	<b>Delivered frame rate (Note 1)</b>	<b>Typical latency Delay variation (Note 2)</b>	<b>Nominal video bit rate (Mbit/s)</b>
TV 0	ITU-R Rec. BT.601	Max FR	(Note 2)	270
TV 1	ITU-R Rec. BT.601	Max FR	(Note 2)	18 to 50
TV 2	ITU-R Rec. BT.601	Max FR	(Note 2)	10 to 25
TV 3	ITU-R Rec. BT.601	Max FR occasional Frame repeat	(Note 2)	1.5 to 8
MM 4a	ITU-R Rec. BT.601	~30 or ~25 fps	Delay $\leq$ 150 ms Variation $\leq$ 50 ms	~1.5
MM 4b	CIF	~30 or ~25 fps	Delay $\leq$ 150 ms Variation $\leq$ 50 ms	~0.7
MM 5a	CIF	10-30 fps	Delay $\leq$ 1000 ms Variation $\leq$ 500 ms	~0.2
MM 5b	$\leq$ CIF	1-15 fps	Delay $\leq$ 1000 ms Variation $\leq$ 500 ms	~0.05
MM 6	CIF-16CIF	Limit $\rightarrow$ 0 fps	No restrictions	<0.05, Limit $\rightarrow$ 0 fps

NOTE 1 – Normally 30 fps for 525 systems and 25 fps for 625 systems

NOTE 2 – Broadcast systems all have constant, but not necessarily low, one-way latency and constant delay variation. For most broadcast applications latency will be low, say between 50 and 500 ms for high quality videoconferencing, and conversational types of applications in general, latency should be preferably less than 150 ms (see Recommendation G.114). Delay variations are allowed within the given range but should not lead to perceptually disturbing time-warping effects.

## APPENDIX I

### Bibliography

- [I.1] GONZALEZ (R.C.), WINTZ (P.), Digital Image Processing, 2nd Edition, *Addison-Wesley Publishing Co.*, Reading, Massachusetts, 1987.
- [I.2] RACE Industrial Consortium Project 1018 HIVITS, WP B5, Picture Quality Measurement, 1988.
- [I.3] Grahm-Field Catalogue Number 13-1240.
- [I.4] Pseudo Isochromatic Plates, engraved and printed by *The Beck Engraving Co., Inc.*, Philadelphia and New York, United States.
- [I.5] KIRK (R.E.), Experimental Design – Procedures for the Behavioural Sciences, 2nd Edition, *Brooks/Cole Publishing Co.*, California, 1982.
- [I.6] VIRTANEN (M.T.), GLEISS (N.), GOLDSTEIN (M.), On the use of Evaluative Category Scales in Telecommunications, HFT 1995, *Human Factors in Telecommunication Conference*, Melbourne, 1995.
- [I.7] GUILFORD (P.), Psychometric methods – McGraw-Hill, New York, 1954.

## APPENDIX II

### Test sequences

The selection of appropriate test sequences is a key point in the planning of subjective assessment. When results of tests, carried out with different groups of observers or in different laboratories, have to be correlated, it is important that a common set of test sequences is available.

A first set of such sequences is described in Table II.1. In this table the following information is given for each sequence:

- the category (defined in Table A.1);
- a brief description of the scene;
- the source format (either 625- or 525-lines, either ITU-R BT.601-4 format or Betacam SP);
- the values of spatial and temporal information (defined in 5.3.1 and 5.3.2 respectively).

All the sequences listed in Table II.1 are in the public domain and may be used freely for evaluations and demonstrations. Some of the sequences suggested belong to the CCIR library described in CCIR Report 1213 [10].

Other sequences of the CCIR library could be suitably used for particular applications like those based on video storage and retrieval.

The set of test sequences is still under study. The set of test sequences listed in Table II.1 can be improved or extended in at least two ways:

- 1) sequences representative of a wider range of applications must be included (e.g. mobile videophone, remote classroom, etc.);
- 2) the source format for every sequence should be the ITU-R BT.601-4 format [3] in both 525- and 625-line versions.

**Table II.1/P.910 – Test sequences for video quality assessment in multimedia applications**

Sequence	Category	Description	Source format	SI	TI
washdc	D	Washington DC map with hand and pencil motion	Betacam SP (525-lines)	130.5	17.0
3inrow	C	Men at table, camera pan	Betacam SP (525-lines)	81.7	30.8
vtc1nw	A	Woman sitting reading news story	Betacam SP (525-lines)	56.2	5.3
susie	A	Young woman on telephone	ITU-R BT.601-4 525-/625-lines	58.7	24.6
flower garden	E	Landscape, camera pan	ITU-R BT.601-4 525-/625-lines	227.0	46.4
smity2	B	Salesman at desk with magazine	Betacam SP (525-lines)	154.5	35.1

## APPENDIX III

### Instructions for viewing tests

The following may be used as the basis for instructions to assessors involved in experiments adopting either ACR, DCR or PC methods.

In addition, the instructions should give information about the approximate test duration, pauses, preliminary trials and other details helpful to the assessors. This information is not included here because it depends on the specific implementation.

#### III.1 ACR

Good morning and thank you for coming.

In this experiment you will see short video sequences on the screen that is in front of you. Each time a sequence is shown, you should judge its quality by using one of the five levels of the following scale.

- 5 Excellent
- 4 Good
- 3 Fair
- 2 Poor
- 1 Bad

Observe carefully the entire video sequence before making your judgement.

#### III.2 DCR

Good morning and thank you for coming.

In this experiment you will see short video sequences on the screen that is in front of you. Each sequence will be presented twice in rapid succession: within each pair only the second sequence is processed. At the end of each paired presentation you should evaluate the impairment of the second sequence with respect to the first one. You will express your judgement by using the following scale:

- 5 Imperceptible
- 4 Perceptible but not annoying
- 3 Slightly annoying
- 2 Annoying
- 1 Very annoying

Observe carefully the entire pair of video sequences before making your judgement.

#### III.3 PC

Good morning and thank you for coming.

In this experiment you will see short video sequences on the screen that is in front of you. Each sequence will be presented twice in rapid succession: each time through a different codec. The order of the sequences and the combination of codecs in the pairs vary in a random way. At the end of each paired presentation you should express your preference by ticking one of the boxes shown below. You will tick box 1 if you prefer the first sequence or box 2 if you prefer the second sequence of the pair

1
---

2
---

Observe carefully the entire pair of video sequences before making your judgement.

## APPENDIX IV

### **The simultaneous double stimulus for a continuous evaluation**

The Simultaneous Double Stimulus for a Continuous Evaluation (SDSCE) is suitable to evaluate the effect of sparse impairments, such as transmission errors, on the fidelity of visual information. This method is derived from the SSCQE method described in [4].

#### **IV.1 Test procedure**

The panel of subjects is watching two sequences contemporaneously: one is the reference, the other one is the test condition. If the format of the sequences is SIF or smaller, the two sequences can be displayed side by side on the same monitor; otherwise, two aligned monitors should be used.

Subjects are requested to check the differences between the two sequences and to judge the fidelity of the video information by moving the slider of a handset-voting device. When the fidelity is perfect, the slider should be at the top of the scale range (coded 100), when the fidelity is null, the slider should be at the bottom of the scale (coded 0).

Subjects are aware of which is the reference and they are requested to express their opinion, while they are viewing the sequences, throughout their whole duration.

#### **IV.2 The training phase**

The training phase is a crucial part of this test method, since subjects could misunderstand their task. Written instructions should be provided to be sure that all the subjects receive exactly the same information. They should include explanation about what the subjects are going to see, what they have to evaluate (i.e. difference in quality) and how they express their opinion. Any question from the subjects should be answered in order to avoid as much as possible any opinion bias from the test administrator.

After the instructions, a demonstration session should be run. In this way subjects are made acquainted both with voting procedures and kind of impairments.

Finally, a mock test should be run, where a number of representative conditions are shown. The sequences should be different from those used in the test and they should be played one after the other without any interruption.

When the mock test is finished, the experimenter should check that in the case of test conditions equal to references the evaluations are close to one hundred; if they are not, he should repeat the explanation and repeat the mock test.

#### **IV.3 Test protocol features**

The following definitions apply to the test protocol description:

- *Video Segment (VS)*: A VS corresponds to one video sequence.
- *Test Condition (TC)*: A TC may be either a specific video process, a transmission condition or both. Each VS should be processed according to at least one TC. In addition, References should be added to the list of TCs, in order to make "reference/reference" pairs to be evaluated.

- *Session (S)*: A session is a series of different pairs VS/TC without separation and arranged in a pseudo-random order. Each session contains at least once all the VS and TC but not necessarily all the VS/TC combinations. All the combinations of VS/TC must be voted by the same number of observers (but not necessarily the same observers).
- *Test Presentation (TP)*: A test presentation is a series of sessions to encompass all the combinations VS/TC.
- *Voting period*: Each observer is asked to vote continuously during a session.

#### **IV.4 Data processing**

Once a test has been carried out, one (or more) data file is (are) available containing all the votes of the different sessions (S) representing the whole vote material of the Test Presentation (TP). A first check of data validity can be done by verifying that each VS/TC pair has been addressed and that an equivalent number of votes has been allocated to each of them.

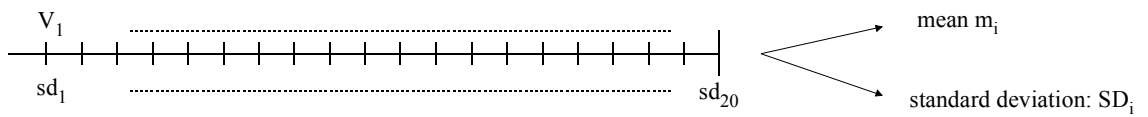
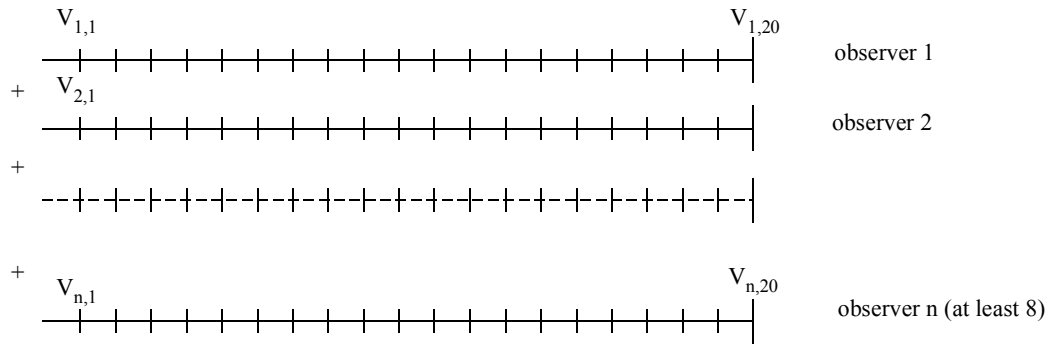
Data of tests carried out according to this protocol can be processed in three different ways:

- Statistical analysis of each separate VS.
- Statistical analysis of each separate TC.
- Overall statistical analysis of all the pairs VS/TC.

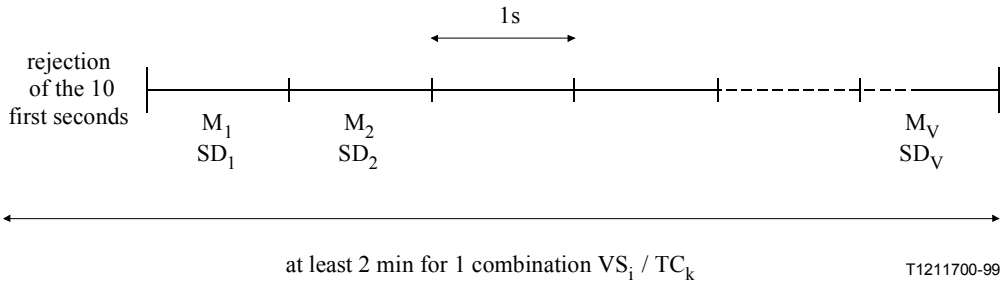
A multi-step analysis is required in each case.

- Means and standard deviations are calculated for each point of vote by accumulation of the observers, as illustrated in Figure IV.1.
- Each VS is then considered as a collection of voting segments of a maximum duration of 10 s. Since neither recency nor forgiveness effect impact the assessment of sequences that lasts no more than 10 s, average and standard deviation of the averages calculated at the previous step are calculated for each voting segment, as illustrated in Figure IV.1. When detailed information about quality variability is required, the duration of voting segment should be short (around one second). The results of this step can be represented in a temporal diagram, as shown in Figure IV.2
- Statistical distribution of the means calculated at the previous step (i.e. corresponding to each voting segment), and their frequency of appearance are analysed. In order to avoid recency effect due to the previous VS/TC, the first 10 seconds of votes for each VS/TC sample are rejected. An example is given in Figure IV.3.
- The global annoyance characteristic is calculated by accumulating the frequencies of occurrence. The confidence intervals should be taken into account in this calculation, as shown in Figure IV.4. A global annoyance characteristic corresponds to this cumulative statistical distribution function by showing the relationship between the means for each voting segments and their cumulative frequency of appearance.

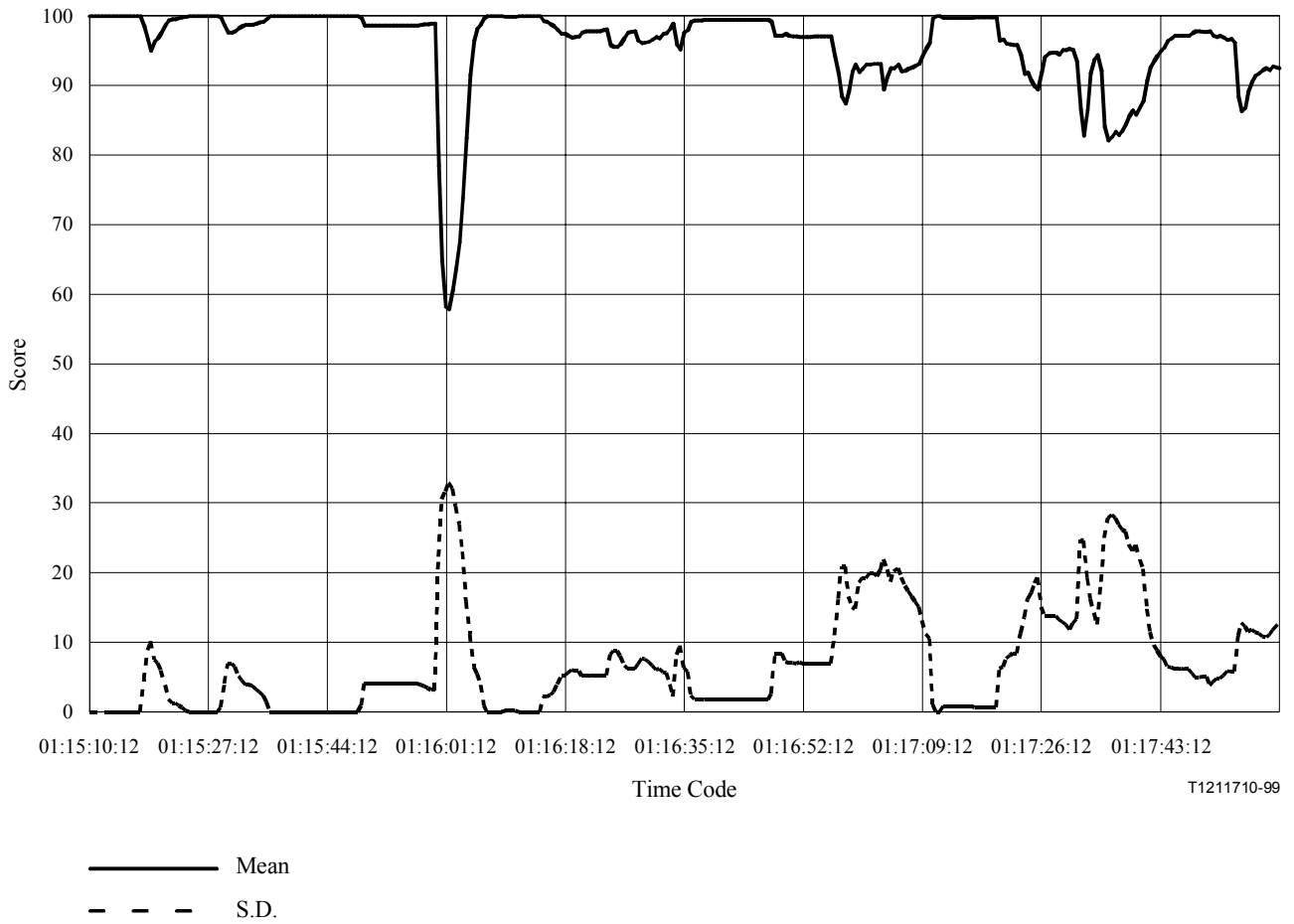
1) Computation of the mean score ( $V$ ) and the standard deviation ( $sd$ ) per instant of vote over observers for every voting sequence of each combination VS/TC.



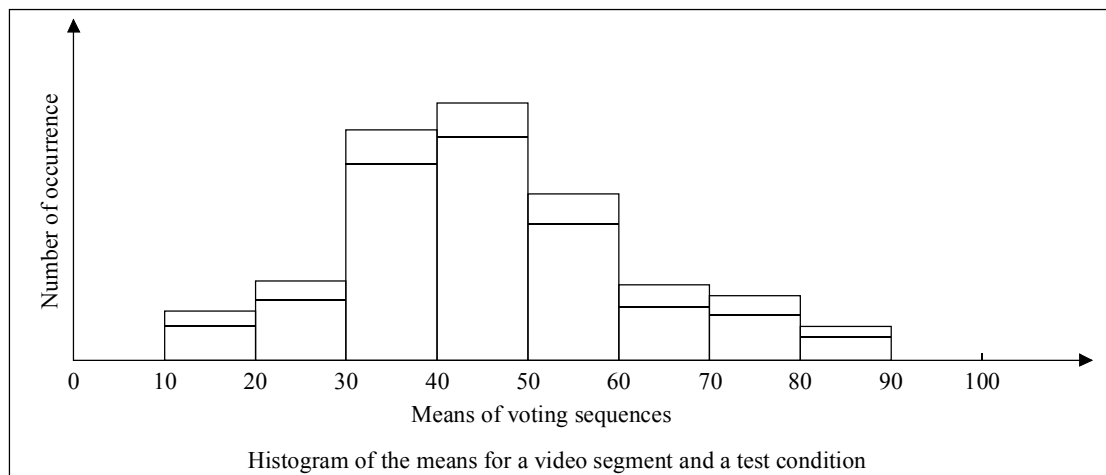
2) Computation of the mean ( $M$ ) and the standard deviation ( $SD$ ) per voting sequence of 1 s for combination VS/TC.



**Figure IV.1/P.910 – Data processing**



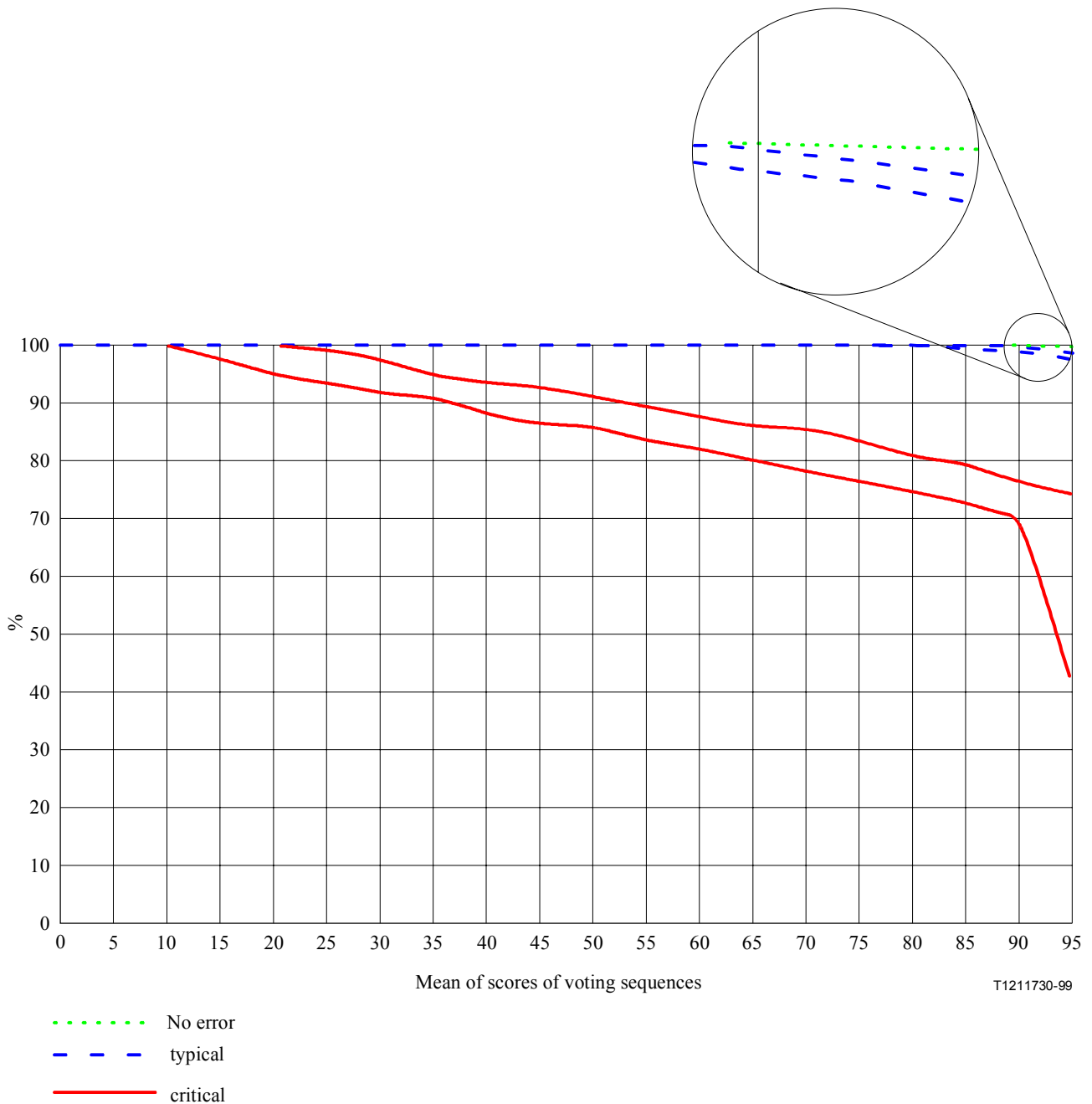
**Figure IV.2/P.910 – Raw temporal diagram**



T1211720-99

**Figure IV.3/P.910 – Relation between the impairment features and their number of occurrence**





**Figure IV.4/P.910 – Global annoyance characteristics calculated from the statistical distributions and including confidence interval**

#### IV.5 Reliability of the subjects

The reliability of the subjects can be qualitatively evaluated by checking their behaviour when "reference/reference" pairs are shown. In these cases, subjects are expected to give evaluations very close to 100. This proves that at least they understood their task and they are not giving random votes.

In addition the reliability of the subjects can be checked by using procedures that are close to that described in [4] for the SSCQE method.

In the SDSCE procedure, reliability of votes depends on the following two parameters:

*Systematic shifts* – During a test, a viewer may be too optimistic or too pessimistic, or may even have misunderstood the voting procedures (e.g. meaning of the voting scale). This can lead to a series of votes systematically more or less shifted from the average series, if not completely out of range.

*Local inversions* – As in other well-known test procedures, observers can sometimes vote without taking too much care in watching and tracking the quality of the sequence displayed. In this case, the overall vote curve can be "relatively" within the average range. But local inversions can nevertheless be observed.

These two undesirable effects (atypical behaviour and inversions) could be avoided. Training of the participants is of course very important. However, the use of a tool allowing to detect and, if necessary, discard inconsistent observers should be possible.

## APPENDIX V

### **The object-based evaluation**

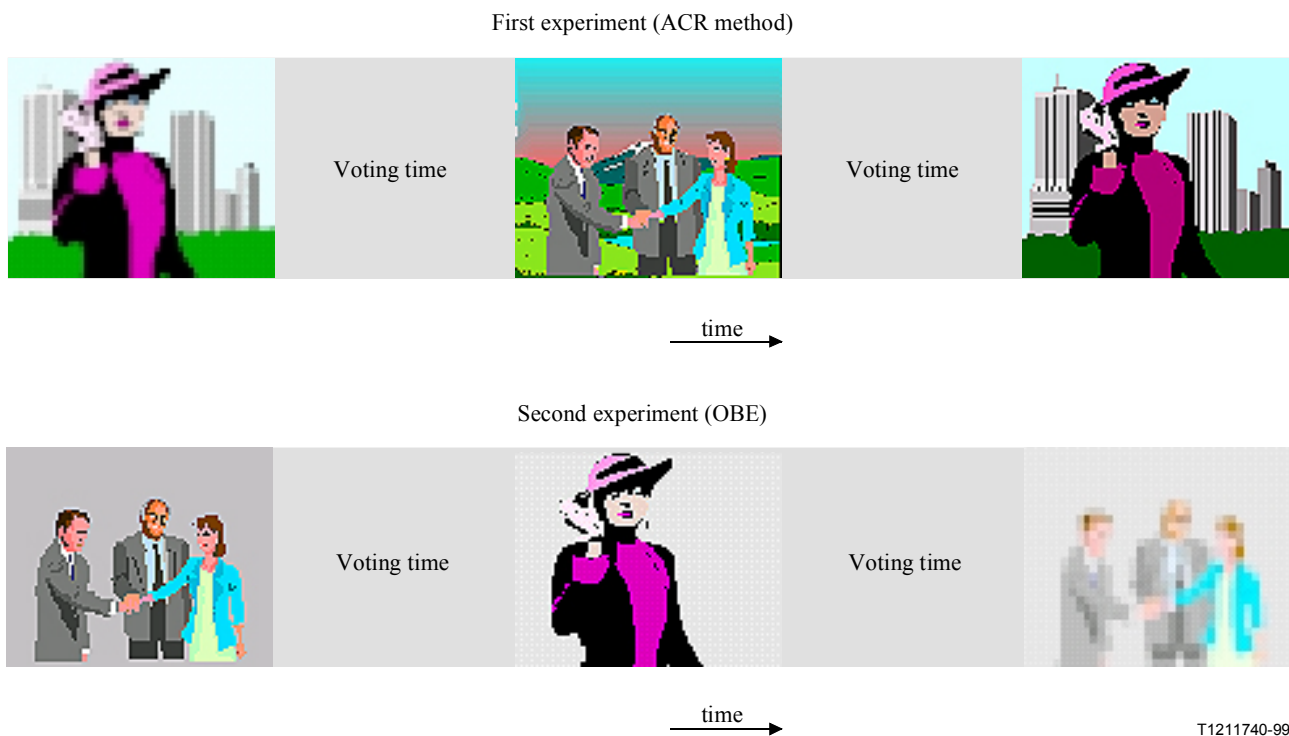
Object-based functionalities should be evaluated both on the whole scene and on the single objects. This is because in general a scene composed of independently encoded objects can be "used" as it has been produced by the author, but in some cases it may also be manipulated and each single object may be used in a completely different context. For this reason it is important to have a balance between the overall quality of the whole scene and the quality of both texture and contours of each single object.

Therefore object-based functionalities (object scalability and object-based quality scalability) should be evaluated in two runs:

*Evaluation of complete picture* – This is a classical test on the whole sequence, that is including all the VOs. The assessment methods may be either the ACR (see 6.1) or the DCR (see 6.2) depending on the range of bit rates and criticality of source sequences.

*Object-Based Evaluation (OBE)* – In this test just one of the VOs will be displayed on a grey background and the subjects will be asked to evaluate the quality/impairment (according to the test method used in the evaluation of the complete picture) of the VO shown. The percentage of bit rate to be spent on the VO has to be specified. The VO evaluated will be extracted from the exact same coded sequence as was used in the complete picture evaluation.

Figure V.1 illustrates the two tests to be carried out for evaluation of object scalability.



**Figure V.1/P.910 – Tests for evaluating object scalability**

In the case of object-based quality scalability, separate tests should be carried out to evaluate spatial scalability and temporal scalability and only OBE should be applied.

Both for spatial and temporal scalability, OBE should be applied to evaluate in the same run both VOs coded at "base" bit rates and the same VOs coded at specified enhanced bit rates.

In general the evaluation of object-based functionalities should take into account both the quality of the whole frame and the quality of the single objects. The former evaluation should be done by standard methods, the latter by means of OBE.

To make a comparison among different systems based on object-based coding, the experimenter should specify in advance the relative weight to assign to global quality and individual object quality.

In particular cases, it will be also worthwhile to use task-based evaluation criteria instead of traditional quality assessments. For example, in the evaluation of a remote monitoring system to be used in a garage, the quality scalability should be evaluated in terms of legibility of car plates. The task will be decided case-by-case by the experimenter, according to the goal of the test and the kind of application under investigation.

Finally, object quality evaluation can be applied to investigate the impact of the quality of the single objects on the overall quality of the scene. Outcomes of such a study could be used to optimize object-based coding schemes.

## APPENDIX VI

### An additional evaluative scale for DRC

A nine-grade degradation scale, as that shown in Figure VI.1, could be used. In this scale grade 8 corresponds to the perceptibility threshold of the degradation, that is the degradation level where the observer is not completely sure to perceive degradation.

9	Imperceptible
8	
7	Perceptible, but not annoying
6	
5	Slightly annoying
4	
3	Annoying
2	
1	Very annoying

**Figure VI.1/P.910 – Nine-grade numerical degradation scale**

## ITU-T RECOMMENDATIONS SERIES

- Series A Organization of the work of the ITU-T
- Series B Means of expression: definitions, symbols, classification
- Series C General telecommunication statistics
- Series D General tariff principles
- Series E Overall network operation, telephone service, service operation and human factors
- Series F Non-telephone telecommunication services
- Series G Transmission systems and media, digital systems and networks
- Series H Audiovisual and multimedia systems
- Series I Integrated services digital network
- Series J Transmission of television, sound programme and other multimedia signals
- Series K Protection against interference
- Series L Construction, installation and protection of cables and other elements of outside plant
- Series M TMN and network maintenance: international transmission systems, telephone circuits, telegraphy, facsimile and leased circuits
- Series N Maintenance: international sound programme and television transmission circuits
- Series O Specifications of measuring equipment
- Series P Telephone transmission quality, telephone installations, local line networks**
- Series Q Switching and signalling
- Series R Telegraph transmission
- Series S Telegraph services terminal equipment
- Series T Terminals for telematic services
- Series U Telegraph switching
- Series V Data communication over the telephone network
- Series X Data networks and open system communications
- Series Y Global information infrastructure and Internet protocol aspects
- Series Z Languages and general software aspects for telecommunication systems