

International Telecommunication Union

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**J.341**

(03/2016)

SERIES J: CABLE NETWORKS AND TRANSMISSION  
OF TELEVISION, SOUND PROGRAMME AND OTHER  
MULTIMEDIA SIGNALS

Measurement of the quality of service – Part 3

---

**Objective perceptual multimedia video quality  
measurement of HDTV for digital cable  
television in the presence of a full reference**

Recommendation ITU-T J.341

ITU-T





## Recommendation ITU-T J.341

### Objective perceptual multimedia video quality measurement of HDTV for digital cable television in the presence of a full reference

#### Summary

Recommendation ITU-T J.341 provides an objective perceptual video quality measurement method for HDTV when a full reference signal is available. The following list shows example applications that can use this Recommendation:

- 1) Potentially real-time, in-service quality monitoring at the source.
- 2) Remote destination quality monitoring when a copy of the source is available at the point of measurement.
- 3) Quality measurement for monitoring of a storage or transmission system that utilizes video compression and decompression techniques, either a single pass or a concatenation of such techniques.
- 4) Lab testing of video systems.

This Recommendation includes an electronic attachment containing several short high definition (HD) video sequences, their corresponding predicted mean opinion scores (MOS), and the predicted MOS for five HD databases available via Video Quality Experts Group (VQEG).

#### History

Edition	Recommendation	Approval	Study Group	Unique ID*
1.0	ITU-T J.341	2011-01-13	9	<a href="http://handle.itu.int/11.1002/1000/11033">11.1002/1000/11033</a>
2.0	ITU-T J.341	2016-03-15	9	<a href="http://handle.itu.int/11.1002/1000/12771">11.1002/1000/12771</a>

---

\* To access the Recommendation, type the URL <http://handle.itu.int/> in the address field of your web browser, followed by the Recommendation's unique ID. For example, <http://handle.itu.int/11.1002/1000/11830-en>.

## FOREWORD

The International Telecommunication Union (ITU) is the United Nations specialized agency in the field of telecommunications, information and communication technologies (ICTs). The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of ITU. 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 Assembly (WTSA), which meets every four years, establishes the topics for study by the ITU-T study groups which, in turn, produce Recommendations on these topics.

The approval of ITU-T Recommendations is covered by the procedure laid down in WTSA Resolution 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.

Compliance with this Recommendation is voluntary. However, the Recommendation may contain certain mandatory provisions (to ensure, e.g., interoperability or applicability) and compliance with the Recommendation is achieved when all of these mandatory provisions are met. The words "shall" or some other obligatory language such as "must" and the negative equivalents are used to express requirements. The use of such words does not suggest that compliance with the Recommendation is required of any party.

## INTELLECTUAL PROPERTY RIGHTS

ITU draws attention to the possibility that the practice or implementation of this Recommendation may involve the use of a claimed Intellectual Property Right. 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, ITU had received notice of intellectual property, protected by patents, which may be required to implement this Recommendation. However, implementers are cautioned that this may not represent the latest information and are therefore strongly urged to consult the TSB patent database at <http://www.itu.int/ITU-T/ipr/>.

© ITU 2016

All rights reserved. No part of this publication may be reproduced, by any means whatsoever, without the prior written permission of ITU.

## Table of Contents

	<b>Page</b>
1 Scope.....	1
1.1 Application .....	2
1.2 Limitations.....	2
2 References.....	2
3 Definitions .....	2
3.1 Terms defined elsewhere .....	2
3.2 Terms defined in this Recommendation.....	3
4 Abbreviations and acronyms .....	3
5 Conventions .....	3
6 Description of the full reference methodology .....	4
7 Findings of the Video Quality Experts Group.....	5
Annex A – Description of model VQuad-HD .....	6
A.1 Preprocessing.....	7
A.2 Time alignment.....	8
A.3 Spatial frame alignment.....	9
A.4 Computation of local similarity and local difference features .....	10
A.5 Analysis of the distribution of local features.....	10
A.6 Computation of the blockiness feature .....	11
A.7 Computation of the jerkiness feature (temporal quality).....	13
A.8 Aggregation to MOS .....	14
A.9 Handling of heavily spatially misaligned video sequences.....	19
A.10 Reference source code implementation.....	20
Annex B – Description of model VQuad-HD-fast.....	21
Annex C – Conformance testing .....	25
Bibliography.....	26



## Recommendation ITU-T J.341

### Objective perceptual multimedia video quality measurement of HDTV for digital cable television in the presence of a full reference

#### 1 Scope

This Recommendation<sup>1</sup> provides a perceptual video quality measurement method for use in high definition television (HDTV) non-interactive applications when the full reference (FR) measurement method can be used. The model VQuad-HD was developed to estimate subjective quality scores obtained using [b-ITU-T P.910].

The full reference measurement method can be used when the unimpaired reference video signal is readily available at the measurement point, as may be the case of measurements on individual equipment or a chain in the laboratory or in a closed environment such as a cable television head end. The estimation method includes both calibration and objective video quality estimations.

The validation test material contained both [b-ITU-T H.264] and MPEG-2 coding degradations and various transmission error conditions (e.g., bit errors, dropped packets). The model proposed in this Recommendation may be used to monitor the quality of deployed networks to ensure their operational readiness. The visual effects of the degradations may include spatial as well as temporal degradations. The model in this Recommendation can also be used for lab testing of video systems. When used to compare different video systems, it is advisable to use a quantitative method (such as that in [b-ITU-T J.149]) to determine the model's accuracy for that particular context.

This Recommendation is deemed appropriate for telecommunication services delivered at between 1 Mbit/s and 30 Mbit/s. The following resolutions and frame rates were considered in the validation test:

- 1080i 60 Hz (29.97 fps);
- 1080p (25 fps);
- 1080i 50 Hz (25 fps);
- 1080p (29.97 fps).

The following conditions were allowed in the validation test for each resolution:

<b>Test factors</b>
Video resolution: 1920 x 1080 interlaced and progressive
Video frame rates 29.97 and 25 frames per second
Video bitrates: 1 to 30 Mbit/s
Temporal frame freezing (pausing with skipping) of maximum 2 seconds
Transmission errors with packet loss
Conversion of the SRC from 1080 to 720p, compression, transmission, decompression, and then conversion back to 1080.
<b>Coding technologies</b>
ITU-T H.264/AVC (MPEG-4 Part 10)
MPEG-2

<sup>1</sup> This Recommendation includes an electronic attachment containing several short high definition (HD) video sequences, their corresponding predicted mean opinion scores (MOS), and the predicted MOS for five HD databases available via VQEG.

Note that 720p was considered in the validation test plan as part of the test condition hypothetical reference circuit (HRC). As 720p is currently commonly up-scaled as part of the display, it was felt that 720p HRCs would more appropriately address this format.

## 1.1 Application

The applications for the estimation model described in this Recommendation include, but are not limited to:

- 1) Potentially real-time, in-service quality monitoring at the source.
- 2) Remote destination quality monitoring when a copy of the source is available at the point of measurement.
- 3) Quality measurement for monitoring of a storage or transmission system that utilizes video compression and decompression techniques, either a single pass or a concatenation of such techniques.
- 4) Lab testing of video systems.

## 1.2 Limitations

The video quality estimation model described in this Recommendation cannot be used to replace subjective testing. Correlation values between two carefully designed and executed subjective tests (i.e., in two different laboratories) normally fall within the range 0.95 to 0.98. If this Recommendation is utilized to make video system comparisons (e.g., comparing two codecs), it is advisable to use a quantitative method (such as that in [b-ITU-T J.149]) to determine the model's accuracy for that particular context.

When frame freezing was present, the test conditions typically had frame freezing durations less than two seconds. The model in this Recommendation was not validated for measuring video quality in a re-buffering condition (i.e., video that has a steadily increasing delay or freezing without skipping). The model was not tested on other frame rates than those used in TV systems (i.e., 29.97 frames per second and 25 frames per second, in interlaced or progressive mode).

It should be noted that in case of new coding and transmission technologies producing artefacts which were not included in this evaluation, the objective model may produce erroneous results. In this case, a subjective evaluation is required.

Note that the model in this Recommendation was not evaluated on talking-head content typical of video-conferencing scenarios.

## 2 References

This Recommendation does not incorporate any normative provisions from external references.

## 3 Definitions

### 3.1 Terms defined elsewhere

This Recommendation uses the following terms defined elsewhere:

**3.1.1 subjective assessment (picture)** [b-ITU-T J.144]: The determination of the quality or impairment of programme-like pictures presented to a panel of human assessors in viewing sessions.

**3.1.2 objective perceptual measurement (picture)** [b-ITU-T J.144]: The measurement of the performance of a programme chain by the use of programme-like pictures and objective (instrumental) measurement methods to obtain an indication that approximates the rating that would be obtained from a subjective assessment test.

**3.1.3 proponent** [b-ITU-T J.144]: An organization or company that proposes a video quality model for validation testing and possible inclusion in an ITU Recommendation.

## **3.2 Terms defined in this Recommendation**

This Recommendation defines the following terms:

**3.2.1 frame rate:** The number of unique frames (i.e., total frames – repeated frames) per second.

**3.2.2 simulated transmission errors:** Errors imposed upon the digital video bit stream in a highly controlled environment. Examples include simulated packet loss rates and simulated bit errors. Parameters used to control simulated transmission errors are well defined.

**3.2.3 transmission errors:** Any errors imposed on the video transmission. Example types of errors include simulated transmission errors and live network conditions.

## **4 Abbreviations and acronyms**

This Recommendation uses the following abbreviations and acronyms:

ACR	Absolute Category Rating
ACR-HR	Absolute Category Rating with Hidden Reference
AVI	Audio Video Interleave
DMOS	Difference Mean Opinion Score
FR	Full Reference
FRTV	Full Reference Television
HD	High Definition
HDTV	High Definition Television
HRC	Hypothetical Reference Circuit
ILG	VQEG's Independent Laboratory Group
MOS	Mean Opinion Score
MOSp	Mean Opinion Score, predicted
NR	No (or zero) Reference
PSNR	Peak Signal-to-Noise Ratio
PVS	Processed Video Sequence
RMSE	Root Mean Square Error
RR	Reduced Reference
SFR	Source Frame Rate
SRC	Source Reference Channel or Circuit
VQEG	Video Quality Experts Group

## **5 Conventions**

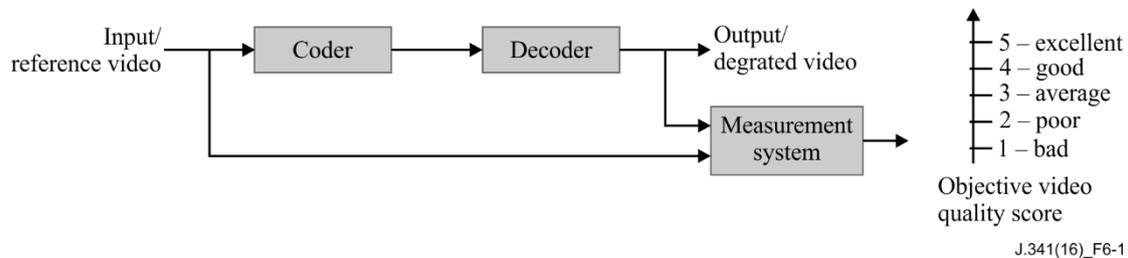
This Recommendation uses the following convention:

YUV      Colour Space and file format

## 6 Description of the full reference methodology

The double-ended measurement method with full reference (FR), for objective measurement of perceptual video quality, evaluates the performance of systems by making a comparison between the undistorted input, or reference, video signal at the input of the system, and the degraded signal at the output of the system (Figure 6-1).

Figure 6-1 shows an example of application of the full reference method to test a codec in the laboratory.



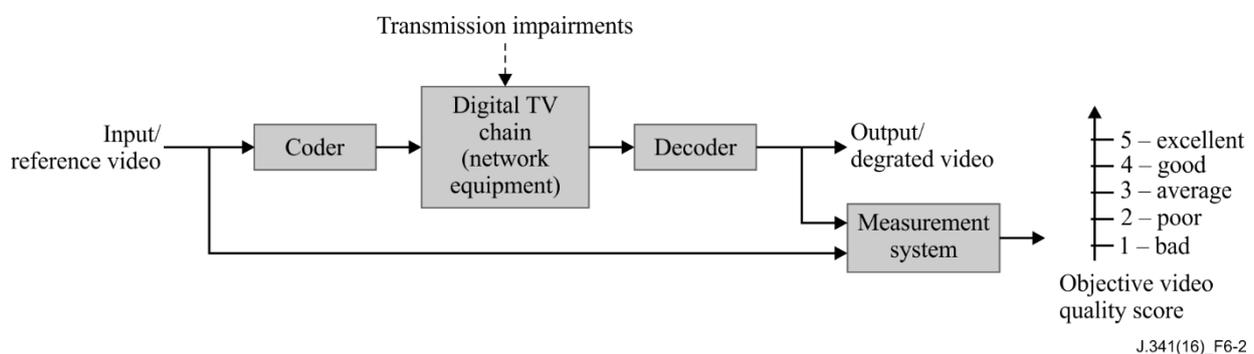
**Figure 6-1 – Application of the full reference perceptual quality measurement method to test a codec in the laboratory**

The comparison between input and output signals may require a temporal alignment or a spatial alignment process, the latter to compensate for any vertical or horizontal picture shifts or cropping. It may also require correction for any offsets or gain differences in both the luminance and the chrominance channels. The objective picture quality rating is then calculated, typically by applying a perceptual model of human vision.

Alignment and gain adjustment is known as registration. This process is required because most full reference methods compare reference and processed pictures on what is effectively a pixel-by-pixel basis. The video quality metrics described in Annexes A and B include registration methods.

As the video quality metrics are typically based on approximations to human visual responses, rather than on the measurement of specific coding artefacts, they are in principle equally valid for analogue and digital systems. They are also in principle valid for chains where analogue and digital systems are mixed, or where digital compression systems are concatenated.

Figure 6-2 shows an example of the application of the full reference method to test a transmission chain.



**Figure 6-2 – Application of the full reference perceptual quality measurement method to test a transmission chain**

In this case, a reference decoder is fed from various points in the transmission chain, e.g., the decoder can be located at a point in the network, as in Figure 6-2, or directly at the output of the encoder, as in Figure 6-1. If the digital transmission chain is transparent, the measurement of objective picture quality rating at the source is equal to the measurement at any subsequent point in the chain.

It is generally accepted that the full reference method provides the best accuracy for perceptual picture quality measurements. The method has been proven to have the potential for high correlation with subjective assessments made in conformity with the absolute category rating with hidden reference (ACR-HR) methods specified in [b-ITU-T P.910].

## 7 Findings of the Video Quality Experts Group

Studies of perceptual video quality measurements are conducted in an informal group, called the Video Quality Experts Group (VQEG), which reports to ITU-T Study Groups 9 and 12 and ITU-R Study Group 6. The recently completed high definition television phase I test of VQEG assessed the performance of proposed full reference perceptual video quality measurement algorithms.

The following statistics are taken from the final VQEG HDTV report [b-VQEG Report]. Note that the body of the VQEG HDTV report includes other metrics including Pearson Correlation and root mean square error (RMSE) calculated on individual experiments, confidence intervals, statistical significance testing on individual experiments, analysis on subsets of the data that include specific impairments (e.g., ITU-T H.264 coding-only), scatter plots, and the fit coefficients.

### Primary analysis

The performance of the FR model is summarized in Table 7-1. The peak signal-to-noise ratio (PSNR) is calculated according to [b-ITU-T J.340] and included in this analysis for comparison purposes. "Superset RMSE" identifies the primary metric RMSE computed on the aggregated superset (i.e., all six experiments mapped onto a single scale). "Top Performing Group Total" identifies the number of experiments (0 to 6) for which this model was either the top performing model or statistically equivalent to the top performing model. "Better Than PSNR Total" identifies the number of experiments (0 to 6) for which the model was statistically better than PSNR. "Better Than Superset PSNR" lists whether each model is statistically better than PSNR on the aggregated superset. "Superset Correlation" identifies the Pearson Correlation computed on the aggregated superset.

**Table 7-1 – Performance of the FR model VQuad-HD**

Metric	PSNR	VQuad-HD
Superset RMSE	0.71	0.56
Top performing group total	1	5
Better than PSNR total	–	4
Better than superset PSNR	–	Yes
Superset correlation	0.78	0.87

## Annex A

### Description of model VQuad-HD

(This annex forms an integral part of this Recommendation.)

This annex contains a description of the model VQuad-HD. Annex B contains a description of the model VQuad-HD-fast, a computationally more efficient variant of VQuad-HD.

#### Overview of the model

The model predicts the video quality as it is perceived by subjects in an experiment. The prediction model uses psycho-visual and cognitive-inspired modelling to emulate subjective perception.

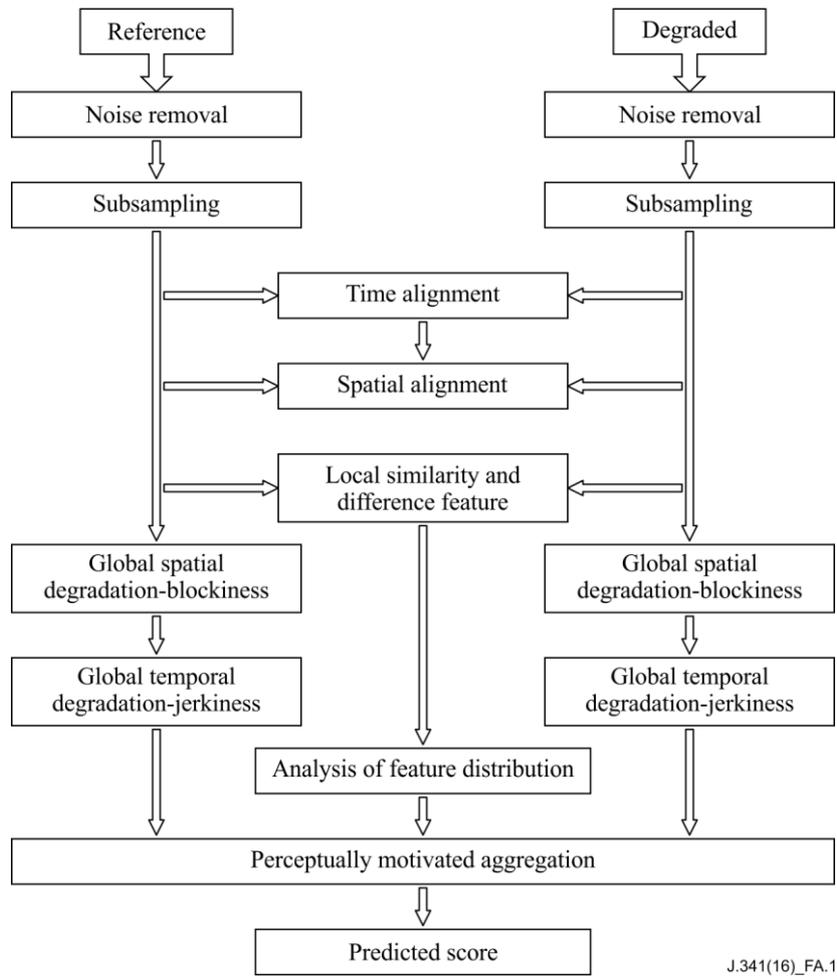
As a full reference approach, the model compares the input or high-quality reference video and the associated degraded video sequence under test. This process is shown in Figure A.1.

Score estimation is based on the following steps:

- 1) First, the video sequences are preprocessed. In particular, noise is removed by filtering the frames and the frames are subsampled.
- 2) A temporal frame alignment between reference and processed video sequence is performed.
- 3) A spatial frame alignment between processed video frame and the corresponding reference video frame is performed.
- 4) Local spatial quality features are computed: a local similarity and a local difference measure, inspired by visual perception.
- 5) An analysis of the distribution of the local similarity and difference feature is performed.
- 6) A global spatial degradation is measured using a blockiness feature.
- 7) A global temporal degradation is measured using a jerkiness feature. The jerkiness measure is computed by evaluating local and global motion intensity and frame display time.
- 8) The quality score is estimated based on a non-linear aggregation of the above features.
- 9) To avoid misprediction in case of relatively large spatial misalignment between reference and processed video sequence (PVS), the above steps are computed for three different horizontal and vertical spatial alignments of the video sequence, and the maximum predicted score among all spatial positions is the final estimated quality score.

The individual steps are explained in more detail in clauses A.1 through A.9. Clause A.10 introduces a C++ source code that covers the essential parts and functions for an implementation-compliant description of the model. The C++ function names mentioned in clauses A.1 through A.9 refer to this reference source code (e.g., clause A.2 refers to `CFrameAnalysisFullRef::ContentTimeAlignment`).

NOTE – On top, the input is the reference and degraded (or processed) video sequences. Different processing steps yield the main model output, the predicted score at the bottom.



**Figure A.1 – Flow chart overview of the processing steps of the model**

### A.1 Preprocessing

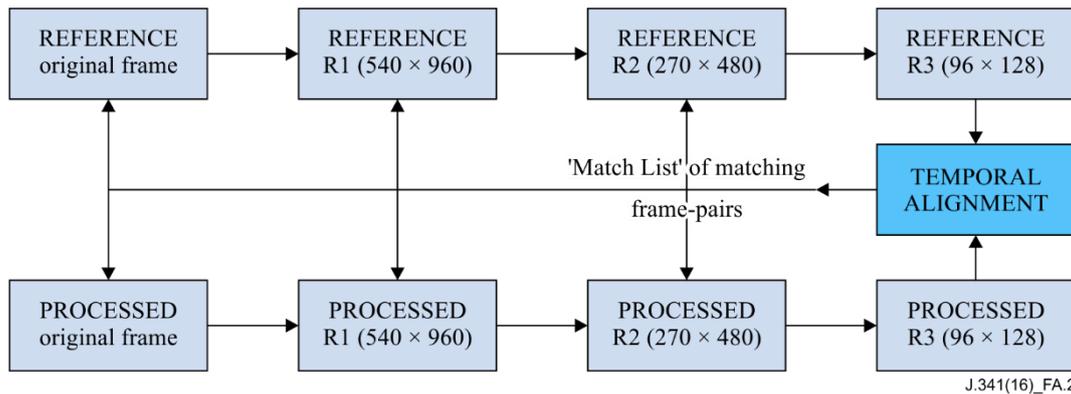
Each frame of the reference and the processed video sequence is spatially lowpass filtered and subsampled to three different resolutions, R1, R2 and R3:

	original frame	R1	R2	R3
height x width	1080 x 1920	→ 540 x 960	→ 270 x 480	→ 96 x 128

See method `CFrameAnalysisFullRef::ContentTimeAlignment` for generating frames at resolution R3 and `CFrameSeq::ReadFrame` for the generation of frames at resolutions R1 and R2.

Note that the implementation is not straightforward due to memory constraints.

Figure A.2 shows the three subsampled resolutions.



NOTE – The frames of the reference and processed video sequence are lowpass filtered and subsampled to three different resolutions. The smallest resolution, R3, is used to perform the frame time alignment. The resulting list of matched frames can be used to match the frames at any other resolution.

**Figure A.2 – Subsampling to obtain frames at three different resolutions**

## A.2 Time alignment

Time alignment is performed using the reference and processed video sequence at the low resolution, R3.

Time alignment is performed in a recursive manner as follows:

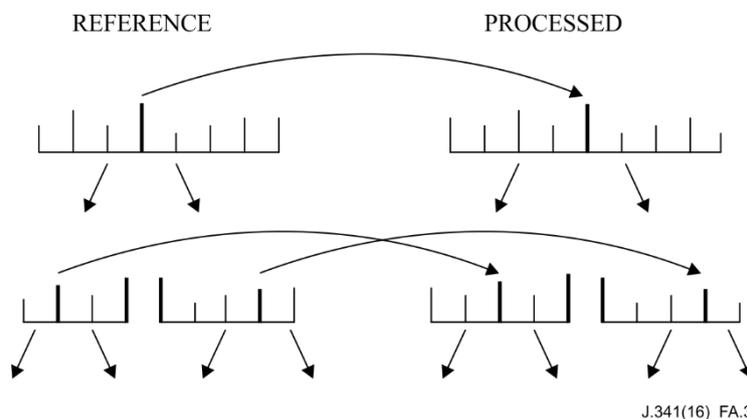
- 1) Find an 'anchor' frame in the reference sequence (Ref\_anchor).
- 2) Match this frame to the best matching frame in the degraded sequence (Deg\_best\_match).

Take this best matched frame in the degraded sequence (Deg\_best\_match) and match it to frames close to the 'anchor' frame of the reference (Ref\_anchor). Try to find a better match, according to a similarity criterion, between the Deg\_best\_match and frames in the environment of Ref\_anchor, and store it as best matching pair. As similarity criterion between the Y-plane of the processed frame  $x$  and the reference frame  $y$ , the function

$$\text{sim} = \exp(-\text{mean\_square\_diff}(a*x+b, y)) \quad (\text{A.2.1})$$

is used, with the parameters  $a$  and  $b$  chosen such that the mean square difference is minimized between the values of the Y-plane of the processed frame  $x$  and the reference frame  $y$ . See method `FrameSimilarity::similarity` in the reference implementation.

- 3) If this matched frame-pair is a good match (the similarity criterion has passed an acceptance threshold), split the reference and processed video sequence at the matching frame-pair, each into two video sequences before and after the matching frames. Start at 1) with both pairs of reference and degraded subsequences.
- 4) If the matching frame-pair is not a 'good match', start again at 1) with a different 'anchor' frame of the reference video. As there is no *a priori* knowledge of the expected value of a 'good' matching frame, the matching threshold is iteratively lowered. The following values were determined based on many training data samples: The starting threshold, with respect to the similarity criterion of equation (A.2.1), is 0.98. After failing to match 10 anchors, it is lowered by a factor of 0.98, and the matching is restarted at 1). This way, at most 10 further anchors are tried, if they fail, the limit is again lowered. This proceeds until reaching a minimum value of 0.1. See `SQ_TimeAlignement::findAnchorAndDescend` for full implementation details.



NOTE – An anchor frame of the reference is matched to a frame of the processed sequence. Then both sequences are split and in each subsequence an anchor frame is chosen and matched.

**Figure A.3 – Recursive approach for time alignment**

The result of the time alignment is a sequence (basically a 'match list'), assigning each frame of the processed video sequence a frame of the reference, or an indicator, indicating that no sufficiently good matching frame could be found. Thus, for the later processing stages, each matched frame of the processed video sequence has a corresponding frame of the reference video. Those frames of the processed video sequence having an indicator 'unmatched' will be compared to the two reference frames matching the previous and following 'matched' frame of the processed video sequence. Note that the 'matching limit' is chosen to be very low, such that only very strongly degraded frames have an indicator 'unmatched'.

See method `CFrameAnalysisFullRef::sqVTA_ContentFrameTimeAlignement_M` for all implementation details.

### A.3 Spatial frame alignment

Iterate over all frames of the processed video sequence and:

- 1) If this frame is unmatched, use the previous spatial alignment. If this frame is matched, perform a spatial alignment between the processed and corresponding – according to the match list of the time alignment – reference frame:
  - a) For the first frame, initialize the spatial shift to be 0 (in both horizontal and vertical directions). For subsequent frames, use as a prior the spatial shift alignment of the previous matched frame.
  - b) Iterate over all possible spatial shifts (horizontal and vertical), using the limit of point 2) below. If a different spatial shift leads to a significantly (with respect to a cost-function) smaller difference between the processed and corresponding reference frame, the spatial shift is adjusted. As a cost function the function

$$\text{rmse}(Y(dv, dh), Y_{\text{ref}}) + \text{abs}(dv) + \text{abs}(dh)$$

is used, where  $Y$  denotes the Y-plane of the processed frame at resolution  $R1$  and  $Y_{\text{ref}}$  denotes the reference frame at resolution  $R1$ ,  $Y(dv, dh)$  denotes the shifted frame  $Y$ , shifted by  $dv$  and  $dh$ , where  $dv$  and  $dh$  are the vertical and horizontal shifts. The second and third terms are included in the cost function, to favour small spatial shifts. Note that a small border of the frames is skipped for the computation of *rmse*, to avoid a more complicated border handling.

- c) That way, time variant spatial shifts can be compensated. Misalignments in one frame can be corrected by the alignment of subsequent frames.

- 2) This first step of the automated spatial shift alignment is limited to  $\pm 4$  pixels. For larger spatial shifts, see clause A.9.
- 3) After the spatial alignment, each frame in the processed video sequence has a corresponding reference frame (or two in the case of an unmatched frame) according to the match-list from time alignment and a well-defined spatial shift correction. Thus, the frames of the processed video sequence can be accurately compared to frames of the reference. This is fundamental for the following feature extractions.

See method `CFrameAnalysisFullRef::DetermineSpatialAlignment` for all implementation details. The constant threshold in step 2 ( $\pm 4$  pixels) can be increased to accommodate larger spatial shifts.

#### A.4 Computation of local similarity and local difference features

For each aligned frame-pair, a set of spatial quality features are calculated:

First, a local similarity and difference measure is computed by iterating over abutting, equally distributed squared regions of size  $13 \times 13$  of the processed and reference frame at resolution R2. As the resolution R2 does not divide by 13, a small border is ignored.

The local regions are called `s_prc` and `s_ref`, and the similarity `S` and difference `D` are computed by:

$$S = (\text{cor}(s\_prc, s\_ref) + 25) / (\text{var}(s\_ref) + 25) \quad (\text{A.4.1})$$

$$D = \text{sqrt}(\text{avg}((S * (s\_prc - \text{mean}(s\_prc)) - (s\_ref - \text{mean}(s\_ref)))^2)) \quad (\text{A.4.2})$$

where `cor` is the correlation and `var` is the variance of the pixel values in the corresponding squared region. The function `avg` computes the average over all pixels of the squared region, and `sqrt` denotes the square root. The values `D` and `S` are the main contributor to the spatial quality value.

At this point, the similarity and difference feature `S`, `D` are a matrix of values for each frame, one value corresponding to each squared local region. Important for the perceived quality is not only the mean value, but the form of the distribution of `S`, `D`, respectively.

#### A.5 Analysis of the distribution of local features

This clause starts with the introduction of some notations:

Let `quantile(X,c)` denote the *c-quantile* of the distribution of the values (entries) of a vector or matrix `X`. More precisely, for a vector `X` and a constant `c` with  $0 \leq c \leq 1$ , the quantile

$$q = \text{quantile}(X, c)$$

is the value `q`, such that a fraction `c` of all the values of `X` is smaller than or equal to `q`.

The function `trimmed_mean` is defined as follows. This notation will be used later. For a matrix `X` the *trimmed\_mean*,

$$\text{trimmedMean}(X, c)$$

is the mean of all entries of `X` between the `c` and  $(1-c)$  quantiles of `X`.

For example, `trimmedMean(X, 0.1)` is the mean of all values of `X`, ignoring 10% of the smallest and 10% of the largest values of `X`.

The notation `X(X>c)` denotes the set of all values of `X` larger than `c`. For example,

$$\text{trimmedMean}(X, c) = \text{mean}(X(X > \text{quantile}(X, c) \text{ and } X < \text{quantile}(X, 1-c)))$$

Using these notations, the following feature values are computed based on `S` from equation (A.4.1), and `D` from equation (A.4.2):

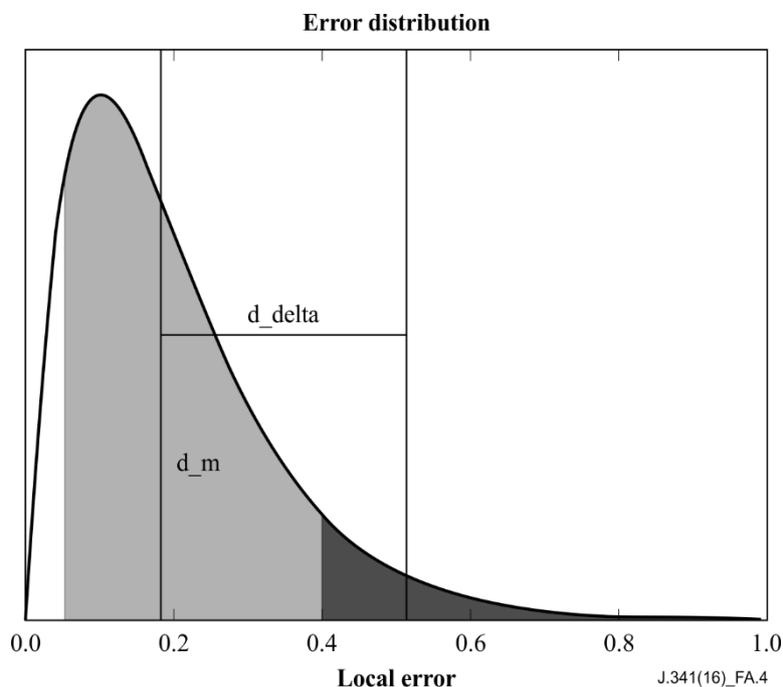
$$s\_m = \text{trimmedMean}(S, c) \quad (\text{A.5.1})$$

$$d\_m = \text{trimmedMean}(D, c) \quad (\text{A.5.2})$$

$$s\_delta = s\_m - \text{mean}(S(S < \text{quantile}(S, c))) \quad (\text{A.5.3})$$

$$d\_delta = \text{mean}(D(D > \text{quantile}(D, 1-c))) - d\_m \quad (\text{A.5.4})$$

using  $c=0.2$ . This is shown visually in Figure A.4, which presents  $d\_m$  and  $d\_delta$ .



NOTE 1 – The trimmed mean  $d\_m$  corresponds to the mean of the light grey region (black vertical line).

NOTE 2 – The value  $d\_delta$  corresponds to the difference of the mean of the values in the dark grey and light grey region (horizontal line).

**Figure A.4 – Distribution of a local D feature**

See method `CFrameAnalysisFullRef::ComputeSimilarity` for the computation of  $s$  and  $D$ .

## A.6 Computation of the blockiness feature

A blockiness feature is computed using the frames at resolution  $R1$ . This feature measures the visibility of block borders introduced by coding and/or transmission errors. Due to the computation at resolution  $R1$ , automatically a focus on the perceptual visibility of edges is set. Starting with an overview, the blockiness feature computes:

- 1) Directional derivatives (edge images) for horizontal and vertical edges. This results in two matrices, one for horizontal and one for vertical edges, for each frame of the video sequence, called `verGrad_n` and `horGrad_n` in the pseudo code below.
- 2) A row-wise and column-wise sum of the logarithm of horizontal and vertical edges, resulting in two vectors, one corresponding to the sum of horizontal, one corresponding to the sum of vertical edges, called `sumW` and `sumH` below.
- 3) An average of a subsample of `sumW` and `sumH`, respectively, at a step size  $n$  and offset  $m$ , computed by the function `vq_AvgSubsample` below.

The idea is that a strong block structure of blocks of size  $n$  will show as an important difference `delta_edge` in the `vq_AvgSubsample` at step size  $n$  computed for different offsets.

For example, a block structure of size 4 in the original frame has a block structure of size 2 at resolution R1. Therefore, computing `vq_AvgSubsample(x, 2, 0)` and `vq_AvgSubsample(x, 2, 1)` should show an important difference, if a strong block structure is present. To avoid content dependency, experiments using a large sample of video sequences showed how to relate the computed difference measured for the processed video sequence to the values of the reference sequence.

More details of the computation are best explained using the following pseudo code. Here, *horGrad* and *verGrad* are the horizontal and vertical spatial derivatives of a frame, given by the difference of adjacent pixels,

$$\text{verGrad}_n(i, j) = Y_n(i+1, j) - Y_n(i, j), \text{ and}$$

$$\text{horGrad}_n(i, j) = Y_n(i, j+1) - Y_n(i, j),$$

where  $Y_n(i, j)$  denotes the pixel value at position  $(i, j)$  of the Y-plane of frame  $n$ . The function

$$\text{vq\_AvgSubsample}(x, \text{step}, \text{offset})$$

calculates the average value of the vector  $x$  over all samples at a step size `step` and starting at offset `offset`.

```
// loop over all frames and compute:
for( UINT i=0; i<horGrad.Height; i++ ){
    for( UINT j=0; j<horGrad.Width; j++ ){
        w = (double)verGrad(i, j);
        h = (double)horGrad(i, j);
        // sum edges (-2: small differences can be the result of integer
        // values used to store frames)
        sumW(i) += log(1.0+max(0.0, fabs(w)-2.0));
        sumH(j) += log(1.0+max(0.0, fabs(h)-2.0));
    }
}

double dH0 = vq_AvgSubsample( sumH, 2, 0 );
double dH1 = vq_AvgSubsample( sumH, 2, 1 );
double dW0 = vq_AvgSubsample( sumW, 2, 0 );
double dW1 = vq_AvgSubsample( sumW, 2, 1 );

edge_max = 0.5 * (vq_Max(dW0, dW1) + vq_Max(dH0, dH1) );
edge_min = 0.5 * (vq_Min(dW0, dW1) + vq_Min(dH0, dH1) );

// Now: denote by edge_max(i) the value of edge_max above, corresponding to
// frame i of the processed video sequence, and by edge_max_ref(i) the values
// edge_max above corresponding to frame i of the reference video sequence,
// and analogously for edge_min(i), edge_min_ref(i). Then compute:

for( UINT i=0; i<nbOfFramesInProcessedVideo; i++ ){
    // get frame nb of ref frame (according to match-list)
    UINT i_ref = (UINT)floor(ref_frameNb_all(i)+0.5f);

    float delta_edge = edge_max(i) - edge_min(i);
    float delta_edge_ref = edge_max_ref(i_ref) - edge_min_ref(i_ref);

    x(i) = vq_Max(0.0f, delta_edge - delta_edge_ref) / (1.0f+edge_max(i));
}
// blockiness(i) is then a non-linear monotone transform of x(i) ...
```

Note that, due to the possibility of upsampled 720-frames, the calculation is slightly more complicated:

See `vquad_hd::vq_BlockinessPhaseDiff` and `CQualityModelFullRef::Blockiness` for all implementation details of the computation of the blockiness feature.

### A.7 Computation of the jerkiness feature (temporal quality)

A jerkiness feature is computed by averaging the product of relative display time, a non-linear transform of display time, and a non-linear transform of motion intensity. The motion intensity is mainly derived by inter-frame differences on individual regions of the frame. *display time* is the time, in milliseconds, a frame is displayed on the screen. Note that to determine the display time of each frame, a local motion intensity analysis is performed, as frames in the processed video sequence might be repetitions of previous frames.

The jerkiness feature takes into account the amount of missed information during playback of the processed video sequence. It is very low in case of a fluently played sequence, while it increases in case of pauses or lowered frame rates. On the other hand, for a fixed temporal impairment, the jerkiness measure takes larger values for video sequences with larger motion intensity.

The following pseudo-code shows the details. Note that the inputs are the motion intensity vector `motionInt`, the vector of frame-repetition probabilities `repFrame` and a vector of frame display times `displayTime`. The output is the vector `jerkiness`, the `jerkiness` at each frame of the processed video sequence. In more detail, the vector `motionInt` denotes the root mean square interframe difference, measured on the Y plane at resolution R2. The vector `repFrame` denotes the probability of frame repetition, i.e., depending on motion intensity, each frame has a probability to be a repetition of the previous frame: in case of a perfect repetition of the previous frame, the actual frame has probability 1 to be a repetition. In case of a large motion intensity, the actual frame has probability 0 to be a repetition of the previous frame. Intermediate values of probabilities can occur if the motion intensity has very small but non zero values. Explicitly,

$$\text{repFrame}(i) = \begin{cases} 0 & \text{if } m(i) < p/2 \\ (m(i)-p/2)/p & \text{if } p/2 \leq m(i) < 3/2*p \\ 1 & \text{if } 3/2*p \leq m(i) \end{cases}$$

where  $m(i)$  denote the motion intensity of frame  $i$ . Empirically, the parameter  $p=0.01$  was chosen.

```
int vq_CalcJerkiness( const CVector<float> & motionInt,
                    const CVector<float> & repFrame,
                    const CVector<float> & displayTime,
                    CVector<float> & jerkiness ){

    // -----
    // the 4 parameters of the jerkiness measure: determined using a
    // large sample of video sequences containing temporal degradations
    // only.
    float a = 0.9f;
    float b = 5.0f;

    float aT = 40.0f;
    float bT = 5.0f;
    // -----

    // get probability of new frame = 1 - prob of repeated frame:
    CVector<float> newFrame = repFrame*(-1.0f) + 1.0f;

    // count number of non-repeated frames
```

```

float fNbRepeated = repFrame.Sum();
float fNbNonRepeated = repFrame.Length() - fNbRepeated;

// calculate jerkiness
float fR = 0.0f;
// look for frame repetition intervals (~ display time)
// of length i
for( UINT i=1; i<= iNbFrames; i++ )
{
    // look for frame repetitions starting at position j
    for( UINT j=0; j<iNbFrames-i+1; j++ )
    {
        float fP = newFrame(j);        // prob. : start of repetition block

        for( UINT k=1; k<i; k++ )
        {
            fP *= repFrame(j+k);      // prob. : all repeated frames
        }
        if( i+j < iNbFrames )
        {
            fP *= newFrame(j+i);      // prob. : end of repetition block
        }

        // calculate the display time (in s) of frame j,
        // if displayed from
        // time t_j until t_(j+i), which occurs with probability fP
        float fDispTime = displayTime.SumPart(j,j+i)/1000.0f;

        // -> measure jerking and add to result
        float fIFDiff = motionInt( j+i-1 );

        // normalisation values: such that at 0 jerkiness value is 0,
        // and saturates at 1
        float c = 1.0f / (1.0f+exp(b));
        float cT = 1.0f / (1.0f+exp(bT));

        float fJ = 1.0f / (1.0f + exp( -( a * fIFDiff - b ) ));
        float fJT = 1.0f / (1.0f + exp( -( aT * fDispTime - bT)));
        fJ = (fJ - c)/(1.0f - c);
        fJT = (fJT - cT)/(1.0f - cT);

        // jerkiness value: propability * interframeDiffFactor *
        //                      displayTimeFactor

        fR = fP * fJ * fJT * fDispTime;

        // add to jerkiness vector at position j+i (corresponding to the
        // end of display time)
        jerkiness( vq_Min(j+i,iNbFrames-1) ) += fR;
    }
}
return 0;
}

```

See method `vquad_hd::vq_ProbOfRepeatedFrame` and `vquad_hd::vq_CalcJerkiness` for the implementation details.

## A.8 Aggregation to MOS

The features described above: Similarity, given by  $s_m$  and  $s_{\Delta}$ , difference, given by  $d_m$  and  $d_{\Delta}$ , blockiness and jerkiness are the basis for score estimation, together with the vector

of frame display times `displayTime`. These are vectors with their lengths equal to the number of frames of the processed video sequence.

To map these features onto a perceptual scale, a parameterized S-shaped function is used:

$$S: \mathbb{R} \rightarrow \mathbb{R}, y = S(x).$$

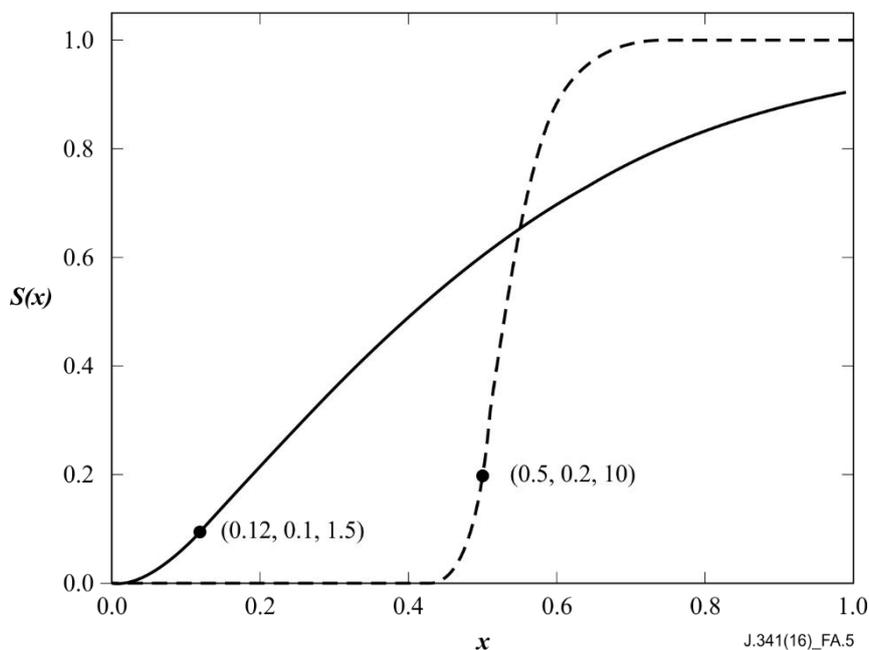
The function  $S$ , is parameterized by the triplet  $(p_x, p_y, q)$ . The parameters have the following interpretation:  $(p_x, p_y)$  is the location in  $\mathbb{R} \times \mathbb{R}$  and  $q$  the slope of the inflection point. In detail:

$$S(x) = \begin{cases} a * x^b & \text{if } x \leq p_x \\ d / (1 + \exp(-c * (x - p_x))) + 1 - d & \text{else} \end{cases} \quad (\text{A.8.1})$$

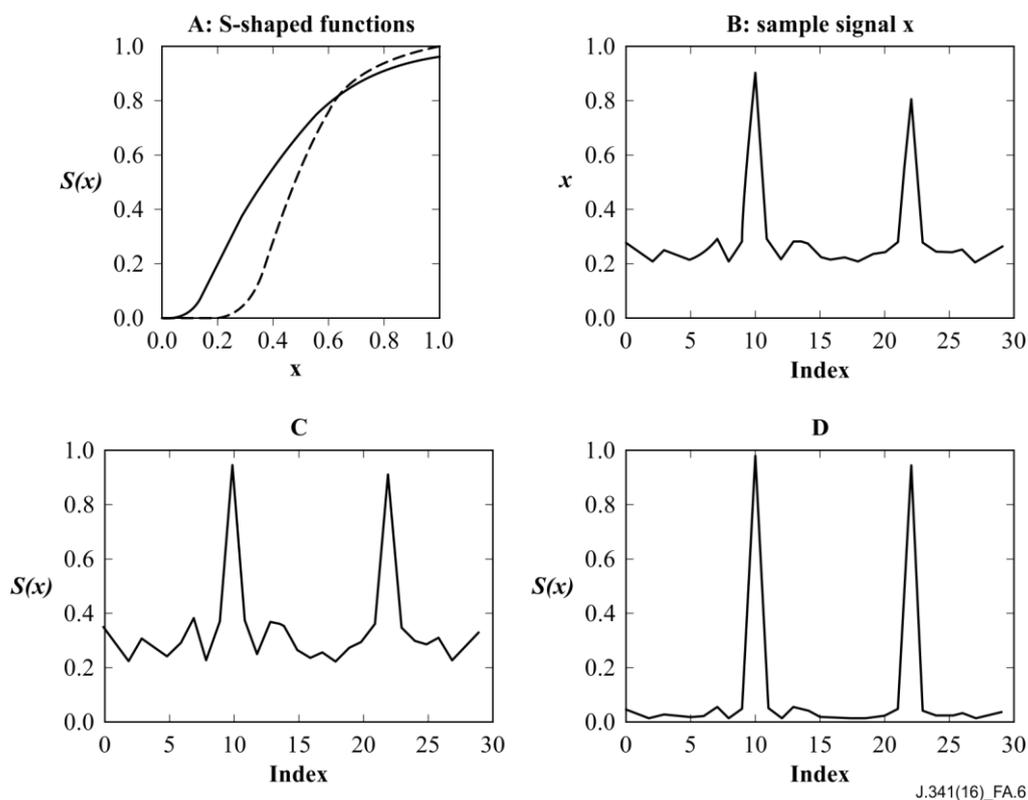
where:

$$\begin{aligned} a &= p_y / p_x^{(q * p_x / p_y)} \\ b &= q * p_x / p_y \\ c &= 4 * q / d \\ d &= 2 * (1 - p_y) \end{aligned}$$

See Figure A.5 for a plot of the S-function with different parameters. The S-shaped function starts at the origin, grows polynomially up to the inflection point and saturates exponentially towards 1.



**Figure A.5 – S-shaped functions, parameterized by the position and slope of the inflection point. Two sample functions are shown for different parameters**



**Figure A.6 – Data-dependent S-shaped function**

Using an S-shaped function with constant parameters, a feature value can be mapped to a perceptual scale. Using data-dependent parameters, the S-shaped function can be used, e.g., to compress all values below the mean to a small value, and stretch all values above. This way, it can be used to compute a feature sensitive to transient degradations, see Figure A.6.

In Figure A.6, plot (A) shows S-shaped functions for two different parameter sets. Plot (B) shows a sample signal vector. All values of  $x$  are transformed using the two S-shaped functions of (A). The result is shown in plots (C) and (D). The S-shaped transformation is used to map a feature  $x$  to a perceptual scale, using parameters fitted to data samples, indicated in (C). Using an S-shaped function (dashed line in (A)) with data dependent parameters (the first parameter is set to a multiple of the signal mean of  $x$ ), a detector for transient degradations can be built, (D).

To continue with the model description, first, define a degradation based on the similarity features above, but putting more weight on the strongest degradations:

$$d_s = 1 - s_m + 1.5 s_{\text{delta}}$$

The next idea is to use two S-shaped functions, one using a fixed set of parameters  $\text{cod\_par}$ , to transform  $d_s$  to  $d_{\text{cod}}$  on a perceptual scale related to a base degradation, reflecting errors due to video encoding.

A second S-shaped function with data-dependent parameters, to transform  $d_s$  to  $d_{\text{trans}}$  on a perceptual scale related to transient degradations, reflecting transmission errors.

In detail, relating to the pseudo-code below, call the function `SplitCodTrans` using as input the vector  $d_s$  and the vector of frame display times  $\text{disp\_time}$  and having output  $d_{\text{cod}}, d_{\text{trans}}$ :

```
SplitCodTrans(d_s, disp_time, d_cod, d_trans)
```

The following pseudo-code shows the function details. Note that

```
stat.STransform(x, px, py, q)
```

denotes the S-shaped function and has as input the real value  $x$  that will be transformed, and the three parameters denoted  $(px, py, q)$  in equation (A.8.1).

```
SplitCodTrans( const CVec& v, const CVec& dispTime,
               CVec& v_cod, CVec& v_trans ){

    // these parameters are determined empirically
    float qPosSmall = 0.55f;
    float qPosLarge = 0.65f;

    // q is the mean of the values in v between the qPosSmall and qPosLarge
    // quantiles.
    float q = r.TrimmedMean( qPosSmall, qPosLarge, v, dispTime );

    for( UINT i=0; i<v.Length(); i++ ){

        // the parameters used here are the result of fitting to sample
        // data
        v_cod(i) = stat.STransform(v(i), 0.07f, 0.1f, 2.0f);

        // Note that the STransform is not directly applied to v, but
        // to v(i)-q . This is preferred here for numerical reasons of the
        // resulting STransform.

        // v_trans is part of v above quantile value q
        v_trans(i) = vq_Max(0.0f, v(i)-q);

        float px = 0.5f * (q+0.2f);
        // the parameters used here are the result of fitting to sample
        // data
        v_trans(i) = stat.STransform(v_trans(i), px, 0.1f, 16.0f);
    }
}
```

See `CQualityModelFullRef::SplitCodTrans` for full implementation details.

In analogy,  $d\_diff\_cod$ ,  $d\_diff\_trans$  are derived from  $d\_m$ ,  $d\_delta$ , by setting

$$d\_diff = d\_m + 1.5 d\_delta,$$

and calling

```
SplitCodTrans(d_diff, disp_time, d_diff_cod, d_diff_trans),
```

using the parameter triples for the two S-transforms

```
cod_par = (4.0f, 0.05f, 0.2f)
trans_par = (0.5*(q+4.0f), 0.1f, 0.4f)
```

where  $q$  denotes the interquantile mean, as in the pseudo-code above. The outputs of the function are  $d\_diff\_cod$ ,  $d\_diff\_trans$ .

Next, a feature value related to transient large jerkiness values is calculated by an S-shaped transform of jerkiness

$$d_{t\_trans} = S(\max(0, \text{jerkiness} - q))$$

with the parameters of the S-shape transform given by  $(\max(0.048, q), 0.2, 40.0)$  and  $q$  denotes the interquantile mean between the 0.55 and 0.65 quantiles of the jerkiness vector. The parameters were determined by fitting to a large set of sample data.

See `CQualityModelFullRef::SplitTempTrans` for full implementation details.

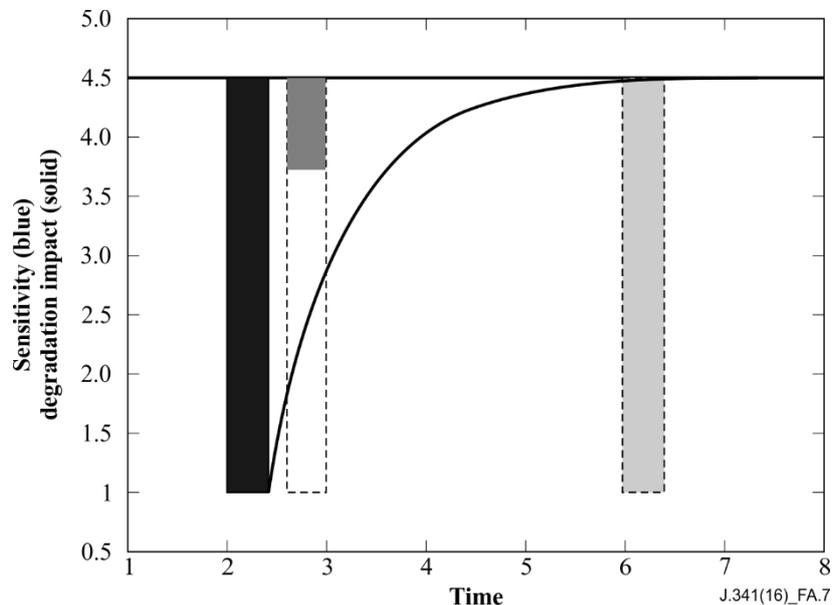
Next, the base quality  $q_{cod}$  is defined as

$$q_{cod} = (1 - d_{cod}) * (1 - d_{diff\_cod}) * (1 - \text{blockiness}) \quad (\text{A.8.2})$$

and the transient quality is defined as

$$q_{trans} = (1 - d_{trans}) * (1 - d_{diff\_trans}) * (1 - d_{t\_trans})$$

The influence of an additional degradation is reduced if it is occurring shortly after a first degradation. To account for this effect,  $q_{trans}$  is transformed to  $q_{fq}$ . The idea is illustrated in Figure A.7.



NOTE – The black solid bar indicates a degradation occurring in a hypothetical video sequence. The black solid line indicates the sensitivity to subsequent degradations. The influence of an additional degradation with strength given by the dashed line is reduced (first dark grey bar) if it occurs shortly after a first degradation. The influence of a subsequent degradation is not altered if occurring with a large time interval in between (right light grey bar).

**Figure A.7 – Reduction of the effect of degradation**

The details of the computation are best described by the following pseudo-code section, which uses  $1 - q_{trans}$  as input vector  $v$ . The vector  $disp\_time$  is the vector of frame display times. The decay time constant  $dT=1000$  ms was determined using sample data; then set

$$q_{fq} = 1 - w \quad (\text{A.8.3})$$

where  $w$  is the output vector of the function:

```
DegFreq( const CVec& v, const CVec& disp_time, CVec& w, float dT ){
    // time constant for temporal integration of degradations
    float t_const = 80.0f;
```

```

w(0) = v(0)*vq_Min(t_const,disp_time(0))/t_const;

for( UINT i=1; i<v.Length(); i++ ){
    // integrate degradation over the last t_const milliseconds:
    // use an indicator function, which is 1 over the interval
    // [t_i-t_const, t_i] and 0 otherwise
    float dT_sum = 0.0f;
    UINT j=0;
    float v_sum = 0.0f;
    while( dT_sum < t_const && (int)i-(int)j>=0 ){
        // b is the overlap in [0,1] of the display time interval
        // of frame i-j with respect to the integration interval
        //           t_i-t_const      t_i
        //           |-----|
        //           |-----> integration window
        //           |-----> time
        //           |-----> frames
        //           ->| b*dT |<-
        //
        float b = vq_Min(t_const-dT_sum,disp_time(i-j)) / t_const;

        v_sum += v(i-j)*b;
        dT_sum += disp_time(i-j);
        j++;
    }
    // compute the fall-off factor a:
    float a = exp(-disp_time(i-1)/dT);
    // the degradation is now:
    // 1) a linear combination of the fall-off of a previous degradation
    //    and the current (integrated) degradation, or
    // 2) the integrated current degradation (if it is stronger than 1)).
    w(i) = vq_Max(v_sum, a * w(i-1) + (1.0f-a) * v_sum);
}

```

## Averaging process

Let  $\text{disp\_time}(i)$  denote the display time of frame  $i$ , and  $\text{jerkiness}(i)$  the jerkiness value corresponding to frame  $i$ . Recall  $q_{\text{cod}}$  from equation (A.8.2) and  $q_{\text{fq}}$  from equation (A.8.3). Set

$$\begin{aligned}
 Q_t &= 1 - 1/T \sum_i [\text{jerkiness}(i)] \\
 Q_{\text{fq}} &= 1/T \sum_i [q_{\text{fq}}(i) * \text{disp\_time}(i)] \\
 Q_{\text{cod}} &= 1/T \sum_i [q_{\text{cod}}(i) * \text{disp\_time}(i)]
 \end{aligned}$$

where  $T = \sum_i [\text{disp\_time}(i)]$ , and  $\sum_i$  denotes the sum over all indices  $i=0, \dots, \text{number\_of\_frames}$ .

The final predicted score is a product and rescaling to the [1,5] mos range:

$$s = 4 * Q_t * Q_{\text{cod}} * Q_{\text{fq}} + 1$$

See method `CQualityModelFullRef::PredictScoreCodTrans` for details about score prediction.

## A.9 Handling of heavily spatially misaligned video sequences

To avoid misprediction in case of relatively large spatial misalignment between reference and processed video sequences, the above steps are computed for three different horizontal and vertical

spatial alignment steps of the video sequence, and the maximum predicted score among all spatial positions is the final estimated quality score.

A step size of four pixels is used in each direction. That way, a spatial search range of  $\pm 8$  pixels is realized. This covers easily the maximum used spatial shift in the test set ( $\pm 5$  pixels). Since, this enlargement is performed in a high-level function of the model, the alignment range can be easily adapted to either larger shifts or can be reduced (e.g.,  $\pm 4$  pixels) for saving computational resources.

See `vquad_hd::vq_vquad08`.

## A.10 Reference source code implementation

The sample video sequences included in this annex are being made available by the Video Quality Experts Group (VQEG) to assist the research community.

The C++ source code covers the required parts and functions for an implementation compliant description of the model. All links to actual implementations made in this Recommendation refer to this reference source code.

The electronic attachment for Recommendation ITU-T J.341 contains the following directories and files:

Directory	Nested file	Contents
Main	Readme.txt	General instructions
	ITU-T_J.341_Conformance.xls	Objective scores predicted by ITU-T J.341 on the sample video sequences
Ref_Code	Copyright_notice_reference_code.txt	Copyright notice for the reference code
	ITU-T_J.341-201101_model_reference_code.cpp	Reference source code for ITU-T J.341 (01/2011)
Video_Sequences	Copyright_notice_video-seq.txt	Copyright notice for the video sequences
	*.avi	Fourteen sample video sequences
	ITU-T_J.341_Conformance_VQEG.xls	Predicted scores for five public VQEG databases

This Recommendation and its electronic attachment are available for free download here: <http://www.itu.int/rec/T-REC-J.341>.

## Annex B

### Description of model VQuad-HD-fast

(This annex forms an integral part of this Recommendation.)

This annex describes the algorithm VQuad-HD-fast, a computationally more efficient alternative to VQuad-HD described in Annex A. Most parts of the algorithm are the same as in VQuad-HD, and overall the changes in the predicted scores are very small. Therefore, the description in this Annex is reduced to the changes to be applied to VQuad-HD to obtain VQuad-HD-fast.

The main modification is in the spatio-temporal alignment of the test video to its reference. This modification does not change the overall performance of the algorithm.

The details of the changes are given in the next section, followed by an analysis of the small differences in score predictions. It will be seen, that these differences are very small and thus irrelevant from a technical point of view. On the other hand, the reduction of computational complexity is important for almost all applications of the algorithm.

#### Description of changes

In this section, two changes to the algorithm VQuad-HD are described to obtain the algorithm VQuad-HD-fast: First a change to the coarse spatial alignment, and second, a change to the rounding of values in subsampling video frame data.

For the first change, instead of performing the coarse spatial alignment in an outer loop, use only one iteration of coarse spatial alignment in the outer loop, and perform the

Old	<pre>verShift.SetValues(-4,-4,-4, 0,0,0, 4,4,4); horShift.SetValues(-4, 0, 4,-4,0,4,-4,0,4);</pre>
New	<pre>verShift.SetValues(0); horShift.SetValues(0);</pre>

#### Figure B.1 – Use one iteration only for the coarse spatial alignment in the outer loop

coarse spatial alignment using low resolution sequences together with the temporal alignment. Apply to method `vq_vquad08(...)` in the reference implementation shown in Figure B.1 and to the temporal alignment the changes shown in Figure B.2. Note the main differences: The low resolution degraded sequence is generated using a spatial shift  $d_{ver}(d)$ ,  $d_{hor}(d)$  in

```
deg.FromFrameSeq(*frameSeq.pBasicFSeq, dver(d), dhor(d),/*iBorderSize=*/0)
```

and therefore, a loop over possible shifts is introduced:

```
for (int d = 0; d<dver.Length(); d++ ){
```

Figure B.2 shows all necessary code changes. *Use coarse spatial alignment together with temporal alignment, in method `CFrameAnalysisFullRef::ContentTimeAlignment(...)`.* Note that large horizontal or vertical shifts are handled by the coarse alignment.

Old	<pre> deg.FromFrameSeq(*frameSeq.pVTA_FSeq,0,0,/*iBorderSize=*/0); deg.DetermineFrameRepetition(settings.repFrameThreshold); iNbMatched = sqVTA_ContentFrameTimeAlignement_M(ref,deg,&amp;settings,matchedFramesRef,matchedFrames,         matchFlag,similarity); </pre>
New	<pre> {     CVector&lt;UINT&gt; matchedFramesRef;     CVector&lt;UINT&gt; matchedFrames;     CVector&lt;UINT&gt; matchFlag;     CVector&lt;float&gt; max_similarity;     float max_mean_similarity = 0.0f;      CVector&lt;int&gt; dver(5);     dver.Fill(0);     CVector&lt;int&gt; dhor(5);     dhor.Fill(0);     dver(1) = -2; dver(2) = 2; // +- 4 pixels in original frame     dhor(3) = -2; dhor(4) = 2; // +- 4 pixels in original frame      for (int d = 0; d&lt;dver.Length(); d++ ){          deg.FromFrameSeq(*frameSeq.pBasicFSeq, dver(d), dhor(d),/*iBorderSize=*/0);         deg.DetermineFrameRepetition(settings.repFrameThreshold);          // call time alignment         iNbMatched = sqVTA_ContentFrameTimeAlignement_M(ref, deg, &amp;settings, matchedFramesRef,             matchedFrames, matchFlag, similarity);          if (similarity.Mean() &gt; max_mean_similarity){             max_mean_similarity = similarity.Mean();             this-&gt;matchedFramesRef = matchedFramesRef;             this-&gt;matchedFrames = matchedFrames;             this-&gt;matchFlag = matchFlag;             max_similarity = similarity;             iVerticalShift = -dver(d);             iHorizontalShift = -dhor(d);              // if a better shift than for (0,0) was found, assume that is the best one             if (d &gt; 0){                 break;             }         }     }     similarity = max_similarity; } </pre>

**Figure B.2 – The coarse spatial alignment is performed using low resolution sequences together with the temporal alignment**

The search range for spatial fine alignment is increased in method `CFrameAnalysisFullRef::Run(...)`, such that large diagonal shifts are covered, as shown in Figure B.3.

Old	<code>float fMaxShift = 2.0f;</code>
New	<code>float fMaxShift = 3.0f;</code>

**Figure B.3 – Increased search range for spatial fine alignment**

Figure B.4 describes the second change whereby an additional performance increase is achieved by avoiding the rounding of float values in sub-sampled video frames. *In method `CFrameSeq::Init(...)`, set the `round values` parameter to `false`.*

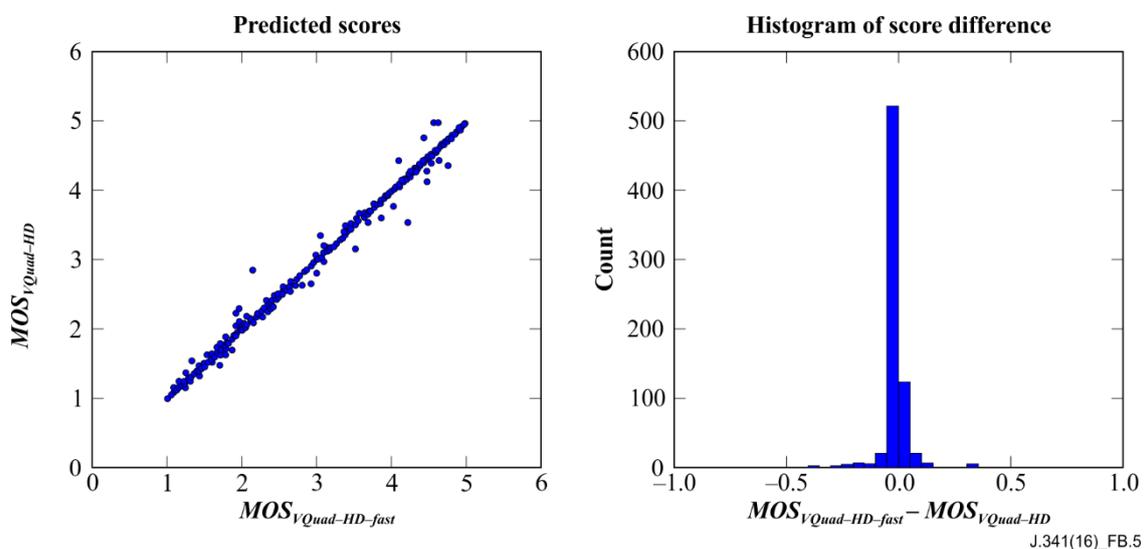
Old	<code>pBasicFSeq-&gt;FromFrameSeq_uchar(*pBasicFSeqOrig,2+iVerticalShift,2+iHorizontalShift,/*iBorderSize=*/0,/*bRoundValues=*/true);</code>
New	<code>pBasicFSeq-&gt;FromFrameSeq_uchar(*pBasicFSeqOrig,2+iVerticalShift,2+iHorizontalShift,/*iBorderSize=*/0,/*bRoundValues=*/false);</code>

**Figure B.4 – Avoid rounding of float values in sub-sampled video frames**

The rounding of frame data was introduced to allow for a smaller memory usage of the algorithm. This advantage is nowadays not necessary anymore, and a faster processing is more desirable.

### Differences in score prediction

The predicted score of the model VQuad-HD is on the MOS scale [1,5] with an RMSE of about 0.56, see Table 7-1. Recommendation ITU-T J.341 contains a conformance test of 714 samples. These samples were evaluated using VQuad-HD as well as VQuad-HD-fast. The Pearson correlation coefficient of the predicted scores of VQuad-HD and VQuad-HD-fast is 0.998. Figure B.5 shows a scatterplot of the predicted values of the two versions and a histogram of the prediction difference. It can be seen, that there is a minor difference between the two versions, keeping in mind that VQuad-HD has an RMSE of about 0.56 on the [1,5] MOS scale.



**Figure B.5 – Predicted score of VQuad-HD versus the predicted score of VQuad-HD-fast and difference in score prediction**

*The scatter plot on the left hand side of Figure B.5 shows the predicted score of VQuad-HD versus the predicted score of VQuad-HD-fast. The histogram at the right of the figure shows the difference in score prediction.*

## Annex C

### Conformance testing

(This annex forms an integral part of this Recommendation.)

This annex is completed by an electronic attachment containing the following information and files:

- 1) Fourteen short HD video sequences (including ten degraded sequences corresponding to four reference sequences). These sequences will cover different distortions and contents. They are provided for testing a conformant model implementation made by a user of this Recommendation against the reference implementation of the model. The video sequences consist of a few frames only to reduce storage capacity. They are not intended for any visual test, but only for a conformance test of the HD model implementation.
- 2) ITU-T\_J.341\_Conformance.xls and ITU-T\_J.341\_VQuad-HD\_fast\_Conformance.xls: File containing predicted MOS scores for HD sequences mentioned under 1). These scores were obtained with a reference implementation of the VQuad-HD model and the VQuad-HD-fast model, respectively.
- 3) ITU-T\_J.341\_Conformance\_VQEG.xls and ITU-T\_VQuad-HD\_fast\_Conformance\_VQEG.xls: File with predicted MOS scores for five public HD databases available via VQEG. These databases can be used as an extended conformance test for implementations of the model. These scores were obtained with a reference implementation of the VQuad-HD model and the VQuad-HD-fast model, respectively.

The attachment can be downloaded from <http://www.itu.int/rec/T-REC-J.341>.

*Conformance test criteria:*

- (I) The 14 scores given in 2) have to be exactly reproduced by an implementation of the model. "Exact" means the reproduction of the score with a resolution of three decimal digits.
- (II) The predicted MOS scores of the five public VQEG databases have to be reproduced with a very limited deviation. Minor variations are allowed, since the experience showed that different optimizations in speed and storage use can lead to minor and negligible deviations in the final score, see table C.1.

**Table C.1 – Allowed distribution of differences across all conformance test data**

Absolute difference	Allowed occurrence
> 0.0001	5.00%
> 0.001	1.00%
> 0.01	0.50%
> 0.1	0.05%
> 0.3	0.00%

For databases different from the ones defined in this Annex C, the same error distribution must not be exceeded. For unknown data, a test set of at least 500 file pairs – preferably from complete subjective experiments – has to be taken for those statistics.

## Bibliography

- [b-ITU-T H.264] Recommendation ITU-T H.264 (2003), *Advanced video coding for generic audiovisual services*.
- [b-ITU-T J.143] Recommendation ITU-T J.143 (2000), *User requirements for objective perceptual video quality measurements in digital cable television*.
- [b-ITU-T J.144] Recommendation ITU-T J.144 (2004), *Objective perceptual video quality measurement techniques for digital cable television in the presence of a full reference*.
- [b-ITU-T J.148] Recommendation ITU-T J.148 (2003), *Requirements for an objective perceptual multimedia quality model*.
- [b-ITU-T J.149] Recommendation ITU-T J.149 (2004), *Method for specifying accuracy and cross-calibration of Video Quality Metrics (VQM)*.
- [b-ITU-T J.244] Recommendation ITU-T J.244 (2008), *Full reference and reduced reference calibration methods for video transmission systems with constant misalignment of spatial and temporal domains with constant gain and offset*.
- [b-ITU-T J.247] Recommendation ITU-T J.247 (2008), *Objective perceptual multimedia video quality measurement in the presence of a full reference*.
- [b-ITU-T J.340] Recommendation ITU-T J.340 (2010), *Reference algorithm for computing peak signal to noise ratio of a processed video sequence with compensation for constant spatial shifts, constant temporal shift, and constant luminance gain and offset*.
- [b-ITU-T P.910] Recommendation ITU-T P.910 (2008), *Subjective video quality assessment methods for multimedia applications*.
- [b-ITU-T P.911] Recommendation ITU-T P.911 (1998), *Subjective audiovisual quality assessment methods for multimedia applications*.
- [b-ITU-T P.931] Recommendation ITU-T P.931 (1998), *Multimedia communications delay, synchronization and frame rate measurement*.
- [b-ITU-R BT.500-11] Recommendation ITU-R BT.500 (2002), *Methodology for the subjective assessment of the quality of television pictures*.
- [b-ITU-R BT.601-6] Recommendation ITU-R BT.601-6 (2007), *Studio encoding parameters of digital television for standard 4:3 and wide screen 16:9 aspect ratios*.
- [b-VQEG Report] VQEG Final Report of HDTV Phase I Validation Test (2010), *Video Quality Experts Group: Report on the Validation of Video Quality Models for High Definition Video Content*.  
<http://www.its.bldrdoc.gov/vqeg/projects/hdtv>



## SERIES OF ITU-T RECOMMENDATIONS

Series A	Organization of the work of ITU-T
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
<b>Series J</b>	<b>Cable networks and transmission of television, sound programme and other multimedia signals</b>
Series K	Protection against interference
Series L	Environment and ICTs, climate change, e-waste, energy efficiency; construction, installation and protection of cables and other elements of outside plant
Series M	Telecommunication management, including TMN and network maintenance
Series N	Maintenance: international sound programme and television transmission circuits
Series O	Specifications of measuring equipment
Series P	Terminals and subjective and objective assessment methods
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, open system communications and security
Series Y	Global information infrastructure, Internet protocol aspects and next-generation networks, Internet of Things and smart cities
Series Z	Languages and general software aspects for telecommunication systems