
**Information technology — MPEG
systems technologies —**

**Part 10:
Carriage of timed metadata metrics of
media in ISO base media file format**

*Technologies de l'information — Technologies des systèmes MPEG —
Partie 10: Transport de métriques de métadonnées de temporisation
de supports au format de fichier de support en base ISO*

IECNORM.COM : Click to view the full PDF of ISO/IEC 23001-10:2020



IECNORM.COM : Click to view the full PDF of ISO/IEC 23001-10:2020



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2020

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Fax: +41 22 749 09 47
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

Page

Foreword	iv
Introduction	v
1 Scope	1
2 Normative references	1
3 Terms, definitions and abbreviated terms	1
3.1 Terms and definitions.....	1
3.2 Abbreviated terms.....	2
4 Carriage of quality metadata	2
4.1 General.....	2
4.2 Quality metadata.....	2
4.2.1 Definition.....	2
4.2.2 Syntax.....	3
4.2.3 Semantics.....	3
4.3 Quality metrics.....	3
4.3.1 Peak signal to noise ratio (PSNR).....	3
4.3.2 SSIM.....	4
4.3.3 MS-SSIM.....	5
4.3.4 VQM.....	7
4.3.5 PEVQ.....	7
4.3.6 MOS.....	8
4.3.7 Frame significance (FSIG).....	8
5 Carriage of green metadata	9
5.1 General.....	9
5.2 Decoder power indication metadata.....	10
5.2.1 Definition.....	10
5.2.2 Syntax.....	10
5.2.3 Semantics.....	10
5.3 Display power reduction metadata.....	10
5.3.1 General.....	10
5.3.2 Display power indication metadata.....	11
5.3.3 Display fine control metadata.....	11
6 Carriage of coordinates	12
6.1 General.....	12
6.2 2D Cartesian coordinates.....	13
6.2.1 2D Cartesian coordinates sample entry.....	13
6.2.2 Syntax.....	13
6.2.3 Semantics.....	13
6.3 2D Cartesian coordinates sample format.....	14
6.3.1 Syntax.....	14
6.3.2 Semantics.....	14
Annex A (informative) Use cases for carriage of ROI coordinates	15
Annex B (normative) Eigen appearance metric matrix specification	17
Bibliography	21

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents) or the IEC list of patent declarations received (see <http://patents.iec.ch>).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 29, *Coding of audio, picture, multimedia and hypermedia information*.

This second edition cancels and replaces the first edition (ISO/IEC 23001-10:2015), which has been technically revised.

The main changes compared to the previous edition are as follows:

- addition of carriage of special information in new [Clause 6](#) and [Annex A](#) with support for encoded regions of interest;
- ISO/IEC 14496-12 and ISO/IEC 23008-2 moved from Bibliography to Clause 2 and other minor editorial changes to align fully with ISO/IEC Directives Part 2.

A list of all parts in the ISO/IEC 23001 series can be found on the ISO website.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html.

Introduction

This document specifies the carriage of timed metadata in files belonging to the family based on ISO/IEC 14496-12. The families of metadata are 'green' metadata (related to energy conservation), quality measurements of the associated media data (related to video quality metrics) and coordinates describing relationship between media data.

IECNORM.COM : Click to view the full PDF of ISO/IEC 23001-10:2020

[IECNORM.COM](https://www.iecnorm.com) : Click to view the full PDF of ISO/IEC 23001-10:2020

Information technology — MPEG systems technologies —

Part 10:

Carriage of timed metadata metrics of media in ISO base media file format

1 Scope

This document defines a storage format for timed metadata. The timed metadata can be associated with other tracks in the ISO base media file format. Timed metadata such as quality and power consumption information and their metrics are defined in this part for carriage in files based on the ISO base media file format (ISO/IEC 14496-12). The timed metadata can be used for multiple purposes including supporting dynamic adaptive streaming.

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 14496-10, *Information technology — Coding of audio-visual objects — Part 10: Advanced video coding*

ISO/IEC 14496-12, *Information technology — Coding of audio-visual objects — Part 12: ISO base media file format*

ISO/IEC 23001-11, *Information technology — MPEG Systems Technologies — Part 11: Energy-Efficient Media Consumption (Green Metadata)*

ISO/IEC 23008-2, *Information technology — High efficiency coding and media delivery in heterogeneous environments — Part 2: High efficiency video coding*

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

ITU-T Recommendation J.247, *Objective perceptual multimedia video quality measurement in the presence of a full reference*

3 Terms, definitions and abbreviated terms

3.1 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 14496-10 and ISO/IEC 23008 apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <http://www.electropedia.org/>

3.2 Abbreviated terms

FSIG	frame significance
MOS	mean opinion score
MSE	mean signal error
MS-SSIM	multi-scale structural similarity index
ROI	region of interest
PEVQ	perceptual evaluation of video quality
PSNR	peak signal to noise ratio
SSIM	structural similarity index
VQM	video quality metric

4 Carriage of quality metadata

4.1 General

If quality metrics are carried in an ISO base media file format, they shall be carried in the metadata tracks within the ISO base media file format in accordance with ISO/IEC 14496-12. Different metric types and corresponding storage formats are identified by their unique code names. This clause defines those quality metrics.

The metadata track is linked to the track it describes by means of a 'cdsc' (content describes) track reference.

Codes not defined in this document are reserved and files shall use only codes defined here.

4.2 Quality metadata

4.2.1 Definition

Sample Entry Type: 'vqme'

Container: Sample Description Box ('stsd')

Mandatory: No

Quantity: 0 or 1

The sample entry for video quality metrics is defined by the `QualityMetricsSampleEntry`.

The quality metrics sample entry shall contain a `QualityMetricsConfigurationBox`, describing metrics that are present in each sample, and the constant field size that is used for the values. The quality metrics are defined in subclause 4.3.

Each sample is an array of quality values, corresponding one for one to the declared metrics. Each value is padded by preceding zero bytes, as needed, to the number of bytes indicated by `field_size_bytes`.

The `codecs` parameter value for this track as defined in RFC 6381^[6] shall be set to 'vqme'. The sub-parameter for the 'vqme' codec is a list of the metrics present in the track as indicated by the metrics code names, joined by "+", e.g., 'vqme.psnr+mssm'.

4.2.2 Syntax

```
aligned(8) class QualityMetricsSampleEntry()
  extends MetadataSampleEntry ('vqme') {
    QualityMetricsConfigurationBox();
  }

aligned(8) class QualityMetricsConfigurationBox
  extends FullBox('vqmC', version=0, 0){
  unsigned int(8) field_size_bytes;
  unsigned int(8) metric_count;
  for (i = 1 ; i <= metric_count ; i++){
    unsigned int(32) metric_code;
  }
}
```

4.2.3 Semantics

`field_size_bytes` indicates the constant size in byte of the value for a metric in each sample.

`metric_count` the number of metrics for quality values in each sample.

`metric_code` is the code name of the metrics in the sample.

4.3 Quality metrics

4.3.1 Peak signal to noise ratio (PSNR)

4.3.1.1 Definition

PSNR for encoded video sequence is defined based on per-picture mean square error (MSE) differences:

$$MSE = \frac{1}{m \cdot n} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [I(i, j) - K(i, j)]^2$$

where

I is the luma plane of the reference $m \times n$ picture;

K is the luma plane of the reconstructed picture;

i, j are indices enumerating all pixel locations.

The picture-level PSNR is defined as:

$$PSNR = 10 \times \log_{10} \left(\frac{MAX_I^2}{MSE} \right)$$

$$PSNR = 20 \times \log_{10} \left(\frac{MAX_I}{\sqrt{MSE}} \right)$$

where $MAX_I = 2^B - 1$ where B is the number of bits per sample in pictures.

PSNR for a given video sequence is computed as an average of all picture-level PSNR values obtained for all pictures in the sequence, i.e., for a sequence with N pictures:

$$PSNR_{sequence} = \frac{1}{N} \sum_{n=0}^{N-1} PSNR_{picture(n)}$$

Only luma component of the video signal is used for PSNR computation.

NOTE 1 This is the traditional metric referred to as PSNR in academic literature and in the context of video compression research.

NOTE 2 In cases when the spatial resolution of the reference pictures and the reconstructed ones do not match, reconstructed pictures are up-sampled to match the spatial resolution of the reference.

NOTE 3 In cases when the pictures of reconstructed video represent only a subset of pictures in the reference video sequence, reconstructed pictures are replicated to produce time-aligned reconstructed pictures for all pictures in the reference sequence.

4.3.1.2 Metric code name

PSNR quality metric values shall be provided as ones under the 'psnr' metric code name.

4.3.1.3 Sample storage format

Each PSNR metric value shall be stored as an unsigned 16-bit integer value.

4.3.1.4 Decoding operation

Given stored 16-bit integer value x, the corresponding PSNR value (in dB) is derived as follows (expressed in floating point):

$$PSNR = (\text{real}) x / 100; \text{ with the exception of } PSNR = \text{infinity for } x=0$$

4.3.2 SSIM

4.3.2.1 Definition

SSIM for encoded video sequence is defined based on SSIM index map obtained for each picture. Per-picture SSIM index map is computed as follows:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

where

- x is the 8×8 window in the reference picture;
- y is the 8×8 window in the reconstructed picture;
- μ_x is the average sample value for pixels in x;
- μ_y is the average sample value for pixels in y;
- σ_x^2 is the variance computed for pixel values in x;
- σ_y^2 is the variance computed for pixel values in y;
- σ_{xy} is the covariance computed for pixel values in x and y.

and where

$$c_1 = (k_1 L)^2, \quad c_2 = (k_2 L)^2$$

are constants computed using

$$k_1 = 0.01, \quad k_2 = 0.03, \quad \text{and} \quad L = 2^B - 1$$

where B is the number of bits per sample in reference video.

This formula is applied using an 8×8 sliding window and producing a map of SSIM index values for all pixel positions within a picture. The overall SSIM index is then computed as the average of index values in the SSIM map.

This formula is applied only on luma components in each picture.

SSIM for video sequence is computed as an average of all picture-level SSIM values obtained for all pictures in the sequence, i.e., for a sequence with N pictures:

$$SSIM_{sequence} = \frac{1}{N} \sum_{n=0}^{N-1} SSIM_{picture(n)}$$

NOTE 1 This is the traditional metric referred to as SSIM in academic literature and in the context of video compression research^[1].

NOTE 2 The nominal range of SSIM index values is [-1..1].

NOTE 3 In cases when the resolution of the reference pictures and the reconstructed ones do not match, reconstructed pictures are up-sampled to match the resolution of the reference.

NOTE 4 In cases when the pictures of reconstructed video represent only a subset of pictures in the reference video sequence, reconstructed pictures are replicated to produce time-aligned reconstructed pictures for all pictures in the reference sequence.

4.3.2.2 Metric code name

SSIM quality metric values shall be provided under the 'ssim' metric code name.

4.3.2.3 Sample storage format

Each SSIM metric value shall be stored as an unsigned 8-bit integer value.

4.3.2.4 Decoding operation

Given stored 8-bit integer value x , the corresponding SSIM value is derived as follows (expressed in floating point):

$$SSIM = (\text{real}) (x - 127) / 128.$$

4.3.3 MS-SSIM

4.3.3.1 Definition

The MS-SSIM calculation procedure is described in [Figure 1](#). Taking the reference and distorted image signals as the input, the system iteratively applies a low-pass filter and downsamples the filtered image by a factor of 2. The original scale is indexed by $j = 1$ and the highest scale is indexed by $j = M$, for $M-1$ levels of iteration. Further details can be found in Reference [\[2\]](#).

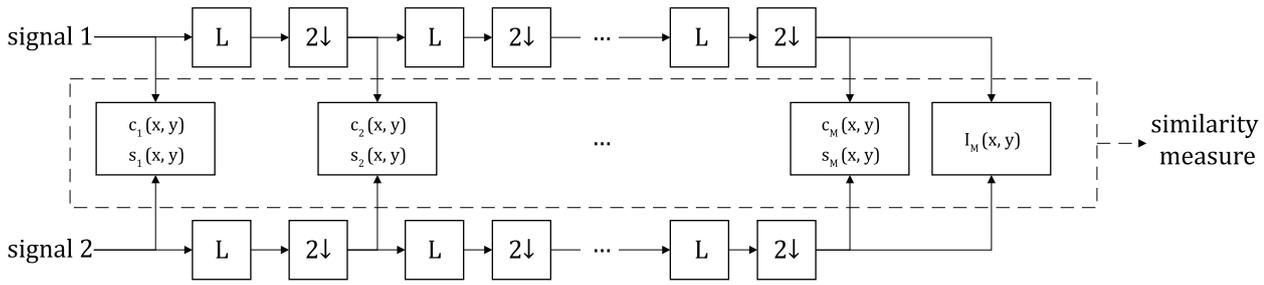


Figure 1 — MS-SSIM calculation procedure

Based on such M scales of processing, MS-SSIM for encoded video sequence is defined as follows:

$$MSSSIM(x, y) = [I_M(x, y)]^{\alpha_M} \prod_{j=1}^M [c_j(x, y)]^{\beta_j} [s_j(x, y)]^{\gamma_j},$$

where

$c_j(x, y)$ is the contrast comparison at scale j ($j = 1, \dots, M$) given by

$$c_j(x, y) = \frac{2\sigma_x \sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2}$$

$s_j(x, y)$ is the structure comparison at scale j ($j = 1, \dots, M$) given by

$$s_j(x, y) = \frac{\sigma_{xy} + C_3}{\sigma_x \sigma_y + C_3}$$

$I_M(x, y)$ is the luma comparison (only computed at scale M) given by

$$I_M(x, y) = \frac{2\mu_x \mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1}$$

where

- x is the 8×8 window in the reference picture;
- y is the 8×8 window in the reconstructed picture;
- μ_x is the average sample value for pixels in x ;
- μ_y is the average sample value for pixels in y ;
- σ_x^2 is the variance computed for pixel values in x ;
- σ_y^2 is the variance computed for pixel values in y ;
- σ_{xy} is the covariance computed for pixel values in x and y .

and where

$$C_1 = (K_1 L)^2, C_2 = (K_2 L)^2, C_3 = C_2/2, \alpha_j = \beta_j = \gamma_j \text{ and } \sum_{j=1}^M \gamma_j = 1$$

are constants computed using

$$k_1 = 0.01, k_2 = 0.03, \text{ and } L = 2^B - 1$$

where B is the number of bits per sample in reference video.

This formula is applied only on luma components in each picture.

MS-SSIM for video sequence is computed as an average of all picture-level MS-SSIM values obtained for all pictures in the sequence, i.e., for a sequence with N pictures:

$$MSSSIM_{sequence} = \frac{1}{N} \sum_{n=0}^{N-1} MSSSIM_{picture(n)}$$

4.3.3.2 Metric code name

MS-SSIM quality metric values shall be provided under the 'msim' metric code name.

4.3.3.3 Sample storage format

Each MS-SSIM metric value shall be stored as an unsigned 8-bit integer value.

4.3.3.4 Decoding operation

Given stored 8-bit integer value x , the corresponding MS-SSIM value shall be derived as follows (expressed in floating point):

$$\text{MS-SSIM} = (\text{real}) (x - 127) / 128$$

4.3.4 VQM

4.3.4.1 Definition

VQM for encoded video sequence is defined as described in ITU-T Recommendation J.144.

4.3.4.2 Metric code name

VQM quality metric values shall be provided under the 'j144' metric code name.

4.3.4.3 Sample storage format

Each VQM metric value shall be stored as an unsigned 8-bit integer value.

4.3.4.4 Decoding operation

Given stored 8-bit integer value x , the corresponding VQM score is derived as follows (expressed in floating point):

$$\text{VQM} = (\text{real}) x / 50$$

4.3.5 PEVQ

4.3.5.1 Definition

PEVQ for encoded video sequence is defined as described in ITU-T Recommendation J.247.

4.3.5.2 Metric code name

PEVQ quality metric values shall be provided as ones carrying 'j247' metric code name.

4.3.5.3 Sample storage format

Each PEVQ metric value shall be stored as an unsigned 8-bit integer value.

4.3.5.4 Decoding operation

Given stored 8-bit integer value x , the corresponding PEVQ score is derived as follows (expressed in floating point):

$$\text{PEVQ} = (\text{real}) x / 50$$

4.3.6 MOS

4.3.6.1 Definition

MOS for encoded video sequence is defined as the arithmetic average of result of a set of standard, subjective tests^[1] where a number of viewers rate the video sequence.

The MOS provides a numerical indication of the perceived quality from the users' perspective of received media after compression. The MOS is expressed as a single number in the range 1 to 5, where 1 is the lowest perceived quality, and 5 is the highest perceived quality. It can be obtained with reference to ITU-R BT.500-12^[6].

4.3.6.2 Metric code name

MOS quality metric values shall be provided as ones under the 'mops' metric code name.

4.3.6.3 Sample storage format

Each MOS metric value shall be stored as an unsigned 8-bit integer value.

4.3.6.4 Decoding operation

Given stored 8-bit integer value x ranging from 0 to 250 (251~255 are reserved), the corresponding MOS value is derived as follows (expressed in floating point):

$$\text{MOS} = \text{ceil}((\text{real}) x / 50)$$

where $\text{ceil}(x)$ is a function which gives the smallest integer not less than x .

4.3.7 Frame significance (FSIG)

4.3.7.1 Definition

FSIG, or frame significance, characterizes the relative importance of frames in a video sequence, and the sequence level visual impact from various combinations of frame losses, e.g., from dropping a temporal layer, can be estimated from this frame significance representation.

For a sequence with frames $\{f_1, f_2, \dots, f_n\}$, The frame significance (FSIG) for frame f_k is defined as:

$$v_k = d(f_k, f_{k-1})$$

where $d()$ is the frame difference function of two successive frames in the sequence.

It is a differential function that captures the rate of change in the sequence^[5], and it is computed from the Eigen appearance metric of the scaled thumbnails^{[4][5]} of the frames:

$$d(f_j, f_k) = (S^* f_j - S^* f_k)^T A^T A (S^* f_j - S^* f_k)$$

where

S is the bi-cubic smoothing and down-scaling function that brings the frames to the size of $h \times w$ pixels;

A is a metric of size $d \times (h \times w)$, where d is the desired dimension of the metric

The metric A is computed from Eigen appearance modelling of thumbnail frames at size $h = 12$, $w = 16$, and $d = 12$, its values provided in [Annex B](#) shall be used.

To characterize the QoE impact of different temporal layers in a video sequence, the visual impact of frame losses are computed from the FSIG in the following fashion. Let the frame loss index be, $L = \{l_1, l_2, \dots, l_n\}$, where $l_k = 1$ if there is a frame loss at time stamp k , and $l_k = 0$, if no frame loss, then the frame losses induced distortion is computed as:

$$D(L) = \frac{1}{n} \sum_{k=1}^n l_k \sum_{j=k}^{p(k)+1} e^{-a(k-j)} v_j$$

where $p(k)$ is the last frame played in the sequence before the loss at frame time k .

An exponentially decaying weight function with kernel size $a=1$ is introduced to model the temporal masking effects for consecutive frame losses.

4.3.7.2 Metric code name

FSIG quality metric values shall be provided as ones carrying 'fsig' metric code name.

4.3.7.3 Sample storage format

Each FSIG metric value is limited to the max value of 255 and shall be stored as an unsigned 8-bit integer value.

4.3.7.4 Decoding operation

Given stored 8-bit unsigned integer value x , the corresponding FSIG value is directly decoded.

5 Carriage of green metadata

5.1 General

If green metadata is carried in an ISO base media file format, it shall be carried in the metadata tracks within the ISO base media file format. Different green metadata types and corresponding storage formats are identified by their unique sample entry codes.

A metadata track carrying green metadata is linked to the track it describes by means of a 'cdsc' (content describes) track reference.

5.2 Decoder power indication metadata

5.2.1 Definition

Sample Entry Type: 'depi'

Container: Sample Description Box ('stsd')

Mandatory: No

Quantity: 0 or 1

The decoder-power indication metadata is defined in ISO/IEC 23001-11. It provides decoder complexity reduction ratios for the media track to which the metadata track refers by means of 'cdsc' reference.

5.2.2 Syntax

The decoder power indication metadata sample entry shall be as follows.

```
class DecoderPowerIndicationMetaDataSetEntry()  
    extends MetaDataSetEntry ('depi') {  
  
}
```

The decoder-power indication sample shall conform to the following syntax:

```
aligned(8) class DecoderPowerIndicationMetaDataSetEntry() {  
    unsigned int(8) Dec_ops_reduction_ratio_from_max;  
    signed int(16) Dec_ops_reduction_ratio_from_prev;  
}
```

5.2.3 Semantics

Semantics are defined in ISO/IEC 23001-11.

5.3 Display power reduction metadata

5.3.1 General

The display-power reduction metadata is defined in ISO/IEC 23001-11. The display power reduction metadata provides frame statistics and quality indicators for the media track that the metadata track refers to by means of 'cdsc' reference. This metadata allows the client to attain a specified quality level by scaling frame-buffer pixels and to reduce power correspondingly by decreasing the display backlight or OLED voltage.

Display-power reduction metadata is of two types:

- metadata that indicates power saving at different quality levels over the sample duration. This metadata shall use the 'dipi' (display power indication) sample entry type.
- metadata that allows fine control of the display to achieve power reduction at a specified quality level. This metadata shall use the 'dfce' (display fine control) sample entry type.

Static metadata for the display fine control is stored in the sample entry. Dynamic metadata is stored in the samples.

5.3.2 Display power indication metadata

5.3.2.1 Definition

Sample Entry Type: 'dipi'

Container: Sample Description Box ('stsd')

Mandatory: No

Quantity: 0 or 1

This metadata indicates potential power saving at different quality levels over the sample duration.

5.3.2.2 Syntax

Display power indication metadata shall use the following sample entry:

```
aligned(8) class DisplayPowerIndicationMetaDataSetEntry() extends MetaDataSetEntry
('dipi') {
}
```

The display power indication sample shall use the following syntax:

```
class QualityLevels (num_quality_levels) {
    unsigned int(8) rgb_component_for_infinite_psnr;
    for (i = 1; i <= num_quality_levels; i++) {
        unsigned int(8) max_rgb_component;
        unsigned int(8) scaled_psnr_rgb;
    }
}
aligned class DisplayPowerIndicationMetaDataSetSample () {
    unsigned int(4) num_quality_levels;
    unsigned int(4) reserved=0;
    QualityLevels(num_quality_levels)
}
```

The PSNR variables appearing in the syntax above are as defined in ISO/IEC 23001-11 and should not be confused with the PSNR metric defined in subclause [4.2](#).

5.3.2.3 Semantics

Semantics are defined in ISO/IEC 23001-11.

5.3.3 Display fine control metadata

5.3.3.1 Definition

Sample Entry Type: 'dfce'

Container: Sample Description Box ('stsd')

Mandatory: No

Quantity: 0 or 1

The display fine control dynamic metadata is stored in the samples and is associated with one or more video frames.

The decoding time to sample box provides the decoding time for the sample so that the metadata contained therein is made available to the display with sufficient lead time relative to the video composition time. Note that the video composition time and metadata composition time are identical. The lead time is required because display settings must be adjusted in advance of presentation time for

correct operation. If `num_constant_backlight_voltage_time_intervals > 1`, then the lead time should be larger than the largest `constant_backlight_voltage_time_interval`.

5.3.3.2 Syntax

The display fine control metadata sample entry shall store static metadata as follows.

```
class DisplayFineControlMetaDataSetEntry()
    extends MetaDataSetEntry ('dfce') {
    DisplayFineControlConfigurationBox();
}

aligned(8) class DisplayFineControlConfigurationBox
    extends FullBox('dfcC', version = 0, flags = 0) {
    unsigned int(2) num_constant_backlight_voltage_time_intervals;
    unsigned int(6) reserved = 0;
    unsigned int(16) constant_backlight_voltage_time_interval[
        num_constant_backlight_voltage_time_intervals ];
    unsigned int(2) num_max_variations;
    unsigned int(6) reserved = 0;
    unsigned int(16) max_variation[ num_max_variations ];
}
```

The display fine control metadata sample shall use the following syntax:

```
class QualityLevels (num_quality_levels) {
    unsigned int(8) rgb_component_for_infinite_psnr;
    for (i = 1; i <= num_quality_levels; i++) {
        unsigned int(8) max_rgb_component;
        unsigned int(8) scaled_psnr_rgb;
    }
}

class MetadataSet (num_quality_levels) {
    unsigned int(8) lower_bound;
    if (lower_bound > 0)
        unsigned int(8) upper_bound;
    QualityLevels(num_quality_levels);
}

class DisplayPowerReductionMetaDataSetEntry
    unsigned int(4) num_quality_levels;
    unsigned int(4) reserved = 0;

    for (k=0; k<num_constant_backlight_voltage_time_intervals; k++)
        for (j = 0; j < num_max_variations; j++)
            MetadataSet(num_quality_levels);
}
```

5.3.3.3 Semantics

Semantics are defined in ISO/IEC 23001-11.

6 Carriage of coordinates

6.1 General

This document specifies the carriage of ROI coordinates in the ISO base media file format using metadata tracks. Different coordinate types and corresponding storage formats are identified by their sample entry. This clause defines those coordinates.

The ROI metadata track shall be linked, via track reference, to the track it describes by means of a 'cdsc' (content describes) track reference, and may be linked to one or more tracks carrying the ROI media content it defines by means of a 'eroi' (encoded region-of-interest) track reference.

The ROI described by a sample in the ROI metadata track indicates the position of the ROI in a reference space. The ROI position is then mapped in the video track with respect to the dimensions documented

by the track header (i.e., on a uniformly sampled grid, possibly upsampled to track header width and height) but before the application of the track (or movie) matrix, if any.

EXAMPLE For a video track whose track width and height (before the application of the track or movie matrix, if any) are `video_width` and `video_height`, referenced by a ROI metadata track, then the ROI coordinates (`roi_x`, `roi_y`) and the ROI size (`roi_width`, `roi_height`) in the video space are given by:

$$\text{roi_x} = \text{top_left_x} \times \frac{\text{video_width}}{\text{reference_width}}$$

$$\text{roi_y} = \text{top_left_y} \times \frac{\text{video_height}}{\text{reference_height}}$$

$$\text{roi_width} = \text{width} \times \frac{\text{video_width}}{\text{reference_width}}$$

$$\text{roi_height} = \text{height} \times \frac{\text{video_height}}{\text{reference_height}}$$

NOTE When the vertical and horizontal scaling factors have different values, the aspect ratio of the ROI in the video and the aspect ratio of the ROI in the ROI metadata track samples are not equal.

Additionally, since a ROI metadata track is a non-visual track, the track width and height of the ROI metadata track should be set to zero as specified in the semantics defined in ISO/IEC 14496-12. Consequently, the flag `Track_size_is_aspect_ratio` of the Track Header Box (see ISO/IEC 14496-12:2012, 8.3.2.3) should not be set since the values of width and height set to zero do not indicate a desired aspect ratio.

6.2 2D Cartesian coordinates

6.2.1 2D Cartesian coordinates sample entry

Sample Entry Type: `'2dcc'`

Container: Sample Description Box (`'stsd'`)

Mandatory: No

Quantity: 0 or 1

The 2D Cartesian coordinates sample entry provides spatial information related to the referenced track expressed in a two-dimension Cartesian coordinate system.

[Annex A](#) provides example usage of the 2D Cartesian coordinates sample entry.

6.2.2 Syntax

The 2D Cartesian coordinates sample entry shall be as follows:

```
aligned(8) class 2DCartesianCoordinatesSampleEntry
    extends MetadataSampleEntry ('2dcc') {
        unsigned int(16)  reference_width;
        unsigned int(16)  reference_height;
    }
```

6.2.3 Semantics

`reference_width` and `reference_height` give respectively the width and height of the reference rectangular space in which all ROI coordinates (`top_left_x`, `top_left_y`, `width` and `height`) are computed. These fields allow associating a ROI metadata track with video tracks of different resolutions but representing the same visual source.

6.3 2D Cartesian coordinates sample format

6.3.1 Syntax

The 2D Cartesian coordinates sample shall conform to the following syntax:

```
aligned(8) class 2DCartesianCoordinatesSample() {
    unsigned int(16) top_left_x;
    unsigned int(16) top_left_y;
    unsigned int(16) width;
    unsigned int(16) height;
    unsigned int(1) interpolate;
    unsigned int(7) reserved;
}
```

Sync samples for ROI metadata tracks are samples for which the `interpolate` value is 0.

6.3.2 Semantics

`top_left_x` and `top_left_y` give respectively the horizontal and vertical coordinate of the top-left corner of the rectangle region associated with the media sample of the referenced track.

`width` and `height` give respectively the width and height of the rectangular region associated with the media sample of the referenced track.

`interpolate` indicates the continuity in time of the successive samples. When true, the application may linearly interpolate values of the ROI coordinates between the previous sample and the current sample. When false, there shall not be any interpolation of values between the previous and the current samples.

NOTE When using interpolation, it is expected that the interpolated samples match the presentation time of the samples in the referenced track. For instance, for each video sample of a video track, one interpolated 2D Cartesian coordinates sample is calculated.

Annex A (informative)

Use cases for carriage of ROI coordinates

A.1 Close-up view (video-to-video)

In this scenario, the content provider offers two videos, namely a wide-angle view and a close-up view. The close-up view generally focuses on particularly interesting parts of the scene, e.g., most popular athletes in sport events. But to ensure a satisfying quality of experience for the end-user, it is important to be able to describe the position of the close-up cam with respect to the wide-angle cam at any point in time of the broadcast. This way the end-user application may seamlessly switch from one video to another provide a smooth experience for the end-user. [Figure A.1](#) illustrates this concept for a live broadcast of cycling races.

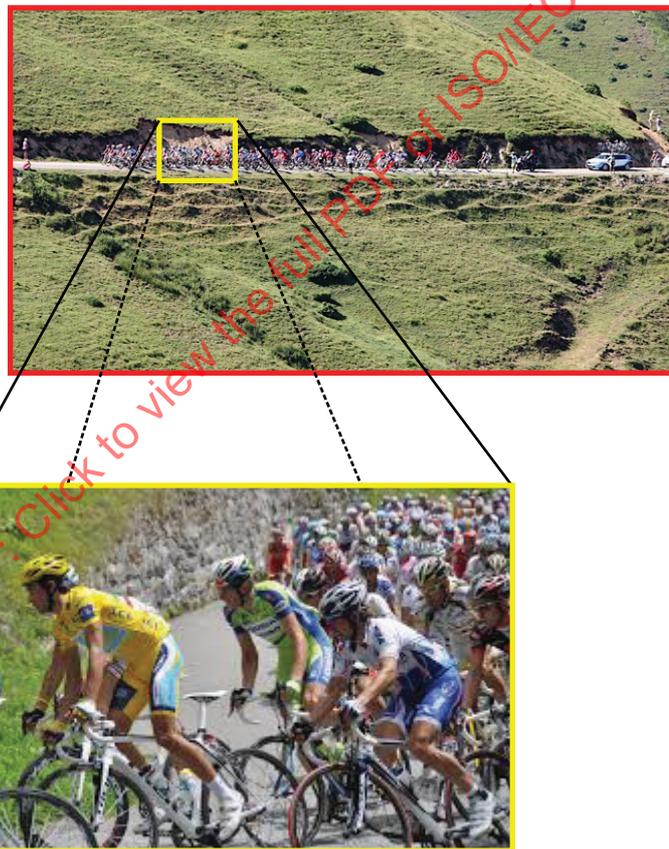


Figure A.1 — "Maillot jaune" cam use case

The file format structure for this scenario is:

- a main video track for the wide-angle view ('main'),
- a second video track offering the close-up view ('closeup'),
- and a timed metadata track ('meta') which contains the time-varying coordinates of the ROI in the 'main' video.

There are then two track references:

- of type `\cdsc` from the 'meta' track to the 'main' video track, saying that the metadata describes the 'main' video track (it describes what the ROI is).
- of type `\eroi` from the 'meta' track to the 'closeup' track, saying that the closeup track is an encoding of the video in only the ROI described by the 'meta' track.

A.2 Object annotation (metadata-to-video)

This scenario involves applications that spatially annotate dynamic object in videos. For instance, a video conference system can provide the position and size of the participants' faces allowing the application to augment the visual experience by displaying participant names, overlaying graphics, etc.

The file format structure for this scenario is the following: a main video offering a view of all the participants, a metadata track containing application specific metadata (e.g., participant information) and timed metadata tracks containing ROI coordinates samples providing the position and size of the participant faces in the main video.

IECNORM.COM : Click to view the full PDF of ISO/IEC 23001-10:2020