

International Telecommunication Union

ITU-T Technical Paper

TELECOMMUNICATION
STANDARDIZATION SECTOR
OF ITU

(3 July 2020)

HSTP-VID-WPOM

**Working practices using objective metrics for
evaluation of video coding efficiency
experiments**

ITU-T

Summary

This ITU-T Technical Paper provides a description of Bjøntegaard Delta rate (BD-rate) measurement practices for video coding experiments. It provides a concept-level overview of recent practices and provides references to technical papers that describe further details. It provides comments on why some of the choices were made and identifies situations where caution must be taken when interpreting the results.

This Technical Paper was developed collaboratively with ISO/IEC JTC 1/SC 29, in technical alignment with ISO/IEC TR 23002-8:2020.

Note

This is an informative ITU-T publication. Mandatory provisions, such as those found in ITU-T Recommendations, are outside the scope of this publication. This publication should only be referenced bibliographically in ITU-T Recommendations.

Keywords

Bjøntegaard Delta, BD-rate, BDR, HEVC, VVC, HDR, 360° video, video coding.

Change Log

This document contains Version 1 of the ITU-T Technical Paper on "*Working practices using objective metrics for evaluation of video coding efficiency experiments*" approved at the ITU-T Study Group 16 virtual meeting held 22 June – 3 July 2020.

Editors:	Jacob Ström Ericsson AB Sweden	E-mail: jacob.strom@ericsson.com
	Kenneth Andersson Ericsson AB Sweden	E-mail: kenneth.r.andersson@ericsson.com
	Rickard Sjöberg Ericsson AB Sweden	E-mail: rickard.sjoberg@ericsson.com
	Andrew Segall Sharp Labs of America USA	E-mail: asegall@sharplabs.com
	Frank Bossen Sharp Labs of America USA	E-mail: fbossen@sharpsec.com
	Gary Sullivan Microsoft USA	E-mail: garysull@microsoft.com
	Jens-Rainer Ohm RWTH Aachen Germany	E-mail: ohm@ient.rwth-aachen.de
	Alexis Tourapis Apple Inc. USA	E-mail: atourapis@apple.com

CONTENTS

	Page
1 SCOPE	1
2 REFERENCES	1
3 TERMS AND DEFINITIONS	2
4 ABBREVIATIONS AND ACRONYMS	2
5 VIDEO CODING EXPERIMENTS USING BJØNTEGAARD DELTA RATE (BD-RATE) MEASUREMENTS	3
6 THE PSNR-BASED BD-RATE CONCEPT	4
7 PSNR-BASED BD-RATE CALCULATION	4
7.1 GENERAL	4
7.2 CALCULATION OF PSNR FOR INDIVIDUAL FRAMES	5
7.3 CALCULATION OF SEQUENCE PSNR AND BIT RATE NUMBERS FOR EACH QP VALUE.....	5
7.4 CALCULATION OF SEQUENCE BD-RATE NUMBER	6
7.5 CONSIDERATION OF CHROMA FIDELITY	8
7.6 CALCULATION OF AGGREGATE BD-RATE VALUE FOR ALL SEQUENCES	9
8 BD-RATE CALCULATION FOR HDR MATERIAL	9
9 BD-RATE CALCULATION FOR 360° VIDEO	10

List of Tables

	Page
TABLE 1 – EXAMPLE BIT RATES AND PSNR VALUES FOR ANCHOR AND TEST	6

List of Figures

	Page
FIGURE 1 – THE LUMA PSNR PLOTTED AS A FUNCTION OF BIT RATE OF HM-16.20 (BLACK) VS VTM-7.0 (RED).....	7
FIGURE 2 – THE LUMA PSNR PLOTTED AS A FUNCTION OF THE BIT RATE OF HM-16.20 (BLACK) VS VTM-7.0 (RED), WHERE THE BIT RATE AXIS IS IN LOG SCALE, AS SUGGESTED IN [3].....	7

Technical Paper ITU-T HSTP-VID-WPOM

ITU-T Technical Paper: Working practices using objective metrics for evaluation of video coding efficiency experiments

1 Scope

This document provides general summary information about coding efficiency measurement practices that have been used for video coding experiments in recent work for the development of video coding Recommendations in ITU-T SG 16. This document does not represent an authoritative recommendation for how video quality should be evaluated. It merely describes the practices that have recently followed for coding efficiency experiments conducted during work to develop video coding standards.

2 References

- [1] G. Bjøntegaard, "Calculation of average PSNR differences between RD-curves," in ITU-T SG 16 Q.6 document VCEG-M33, 13th VCEG meeting, Austin, Texas, USA, Apr. 2001.
- [2] S. Pateux and J. Jung, "An Excel add-in for computing Bjøntegaard metric and its evolution," ITU-T SG 16 Q.6 document VCEG-AE07, 31st VCEG meeting, Marrakech, Morocco, Jan. 2007.
- [3] G. Bjøntegaard, "Improvements of the BD-PSNR model," ITU-T SG16 Q.6 document VCEG-AI11, 35th VCEG meeting, Berlin, Germany, July 2008.
- [4] S. Pateux and J. Jung, "Improvements of Excel macro for BD-gain computation," ITU-T SG16 document C.358, Geneva, Switzerland, Oct. 2009.
- [5] A. M. Tourapis, D. Singer, Y. Su, K. Mammou, "BD-Rate/BD-PSNR Excel extensions," Joint Video Exploration Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 document JVET-H0030, 8th JVET meeting, Macao, China, 18–25 Oct. 2017.
- [6] Y. Ye, J. Boyce, "Algorithm descriptions of projection format conversion and video quality metrics in 360Lib (Version 9)," Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 document JVET-M1004, 13th JVET meeting, Marrakech, Morocco, 9–18 Jan. 2019.
- [7] A. Segall, E. François, W. Husak, S. Iwamura, D. Rusanovskyy, "JVET common test conditions and evaluation procedures for HDR/WCG video," Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 document JVET-P2011, 16th JVET meeting, Geneva, Switzerland, Oct. 2019.
- [8] F. Bossen, "Common HM test conditions and software reference configurations," Joint Collaborative Team on Video Coding (JCT-VC) of ISO/IEC MPEG and ITU-T VCEG document JCTVC-G1200, 7th JCT-VC meeting, Geneva, Switzerland, Nov. 2011.
- [9] F. Bossen, X. Li, A. Norkin, K. Sühring, "JVET AHG report: Test model software development (AHG3)," Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 document JVET-O0003, 15th JVET meeting, Gothenburg, Sweden, July 2019.
- [10] SMPTE ST 2084, *High Dynamic Range Electro-Optical Transfer Function of Mastering Reference Displays*, 2014.
- [11] Recommendation ITU-R BT.2100-2 (2018), *Image parameter values for high dynamic range television for use in production and international programme exchange*.

- [12] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity", *IEEE Transactions on Image Processing*, Vol. 13, No. 4, pp. 600–612, April 2004.
- [13] Z. Wang, E. P. Simoncelli, A. C. Bovik, "Multiscale structural similarity for image quality assessment" in the *Proceedings of the Thirty-Seventh Asilomar Conference on Signals, Systems and Computers*, Vol. 2, pp. 1398–1402, 2004.
- [14] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, M. Manohara, "Toward A Practical Perceptual Video Quality Metric," *Netflix TechBlog*, June 2016.
- [15] J. Ström, K. Andersson, R. Sjöberg, F. Bossen, A. Segall, G. J. Sullivan, J.-R. Ohm, "Suggested content for summary information on BD-rate experiment evaluation practices," Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 input document JVET-Q0826, 17th JVET meeting, Brussels, January 2020.
- [16] K. Andersson, F. Bossen, J.-R. Ohm, A. Segall, R. Sjöberg, J. Ström, G. J. Sullivan, A. Tourapis "Summary information on BD-rate experiment evaluation practices," Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11 output document JVET-R2016, 18th JVET meeting by teleconference, April 2020.

3 Terms and definitions

3.1 Bjøntegaard Delta rate: Average percentage bit rate difference at equal measured distortion, integrated across a range of bit rates in the log domain (as originally described in [1]).

4 Abbreviations and acronyms

This Technical Paper uses the following abbreviations and acronyms:

AVC	Advanced Video Coding (Rec. ITU-T H.264 ISO/IEC 14496-10)
BD-rate	Bjøntegaard Delta rate
HDR	High Dynamic Range
HEVC	High Efficiency Video Coding (Rec. ITU-T H.265 ISO/IEC 23008-2)
HLG	Hybrid Log Gamma
JCT-VC	Joint Collaborative Team on Video Coding (for development of HEVC)
JVET	Joint Video Experts Team (for development of VVC)
MPEG	Moving Picture Experts Group
MS-SSIM	Multi-scale Structural Similarity
MSE	Mean Square Error
PQ	Perceptual Quantizer (as defined in [10] and [11])
PSNR	Peak Signal-to-Noise Ratio
QP	Quantization Parameter
SDR	Standard Dynamic Range
SSIM	Structural Similarity
VMAF	Video Multimethod Assessment Fusion
VVC	Versatile Video Coding
WCG	Wide Colour Gamut
WVGA	Wide Video Graphics Array

$Y'CBCR$	Colour space representation commonly used for video/image distribution, also written as YUV
YUV	Colour space representation commonly used for video/image distribution, also written as $Y'CBCR$

5 Video coding experiments using Bjøntegaard Delta rate (BD-rate) measurements

This document provides general summary information about coding efficiency measurement practices that have been used for video coding experiments in recent work for the development of video coding Recommendations in ITU-T SG 16. Such work has often been conducted together with ISO/IEC JTC 1/SC 29/WG 11 (MPEG) in the JVET and JCT-VC joint collaborative teams. In particular, the document describes the use of Bjøntegaard Delta rate (BD-rate) measurements. It aims to provide a concept-level overview of such recent practices and to provide some references to other works that describe further details. It provides comments on why some of the choices were made and points at situations where caution must be taken when interpreting the results.

This document does not represent an authoritative recommendation for how video quality should be evaluated. It merely describes the practices that have recently been followed for coding efficiency experiments conducted during work to develop video coding standards.

For comparing different encodings, often it is helpful to control the encodings so that similar types and degrees of encoder optimization are applied, except for the aspects to be tested.

When there are large differences between the coding technologies being tested, and especially when there may be a substantial difference between the resulting subjective quality, subjective testing (i.e., using humans to measure the visual quality) is the appropriate action. There are also cases where the quality difference is expected to be primarily a matter of subjective effect – for example, when measuring the effects of deblocking filters.

In the video coding community, we have typically used formal subjective testing at the Call for Proposals and Verification Testing stages of projects (i.e., at the beginning and the end of the work). For measuring smaller effects and where formal subjective testing is not feasible, it is necessary to use objective measurements. Since objective measurements are collected at multiple operational points, and to better understand coding behaviour across all these points, what has commonly been used in our community is the technique known as the BD-rate (Bjøntegaard Delta bit rate) comparison [1].

We ordinarily perform encoding in the $Y'CBCR$ domain (nicknamed YUV herein for brevity and ease of typing). For typical multimedia applications, it is well known that the human visual system is most sensitive to the fidelity of the Y component. The Y component also tends to use most of the bit rate, so it is natural to focus primarily on the Y component. However, we typically measure and report the fidelity of all three components and review the balance between luma and chroma fidelity when interpreting the results. This is also done so as to avoid situations where luma gain may be achieved at a significant cost of chroma fidelity.

To calculate the Bjøntegaard Delta bit rate, a distortion metric needs to be used. For standard-dynamic range video, the distortion metric primarily used in our community has been the Peak Signal to Noise Ratio (PSNR). There are certainly some weaknesses to the PSNR-based BD-rate measure in terms of its correspondence with human perception of fidelity. Some other objective distortion metrics, which claim to have a better relationship with human perception, such as the Structural SIMilarity (SSIM) Index [12], the Multi-Scale SSIM (MS-SSIM) [13], and Video Multimethod Assessment Fusion (VMAF) [14], have also been considered with the BD-rate measurement process. However, the use of PSNR-based BD-rate measures is the most prevalent in the video coding standardization community, for several reasons which we will not try to expound upon here in the interest of brevity and since this document is intended merely to describe the

common practice. In the following we will use the shorter term, BD-rate, to denote PSNR-based BD-rate unless explicitly mentioning another distortion metric.

This document is based on JVET input document [15] and JVET output document [16].

6 The PSNR-based BD-rate concept

When developing a video coding standard, it is important to have a uniform way of reporting the compression results so that different contributions can be compared against each other.

The PSNR metric is based on the squared error of individual sample values and does not take into account how the human visual system works. A relevant question is therefore whether the PSNR metric is a good predictor of subjective quality. The answer depends at least partly on how different the encoding methods being compared are to each other. If the two methods differ greatly, their artefacts may be very different, and the perceived subjective quality will depend heavily on which type of artefact is psychovisually more disturbing. BD-rate measurements are most often used to compare between two versions of the same video encoder that only differ in that in one of them one tool has been turned on or has been modified versus the other. In this scenario it is much more likely that the BD-rate score between these two versions will correlate with a difference in subjective quality. A clear exception is when tools are considered that are only (or primarily) expected to affect subjective quality, such as deblocking filters. Here, decisions are almost always based on a subjective test or expert viewing, and BD-rate numbers are provided more as an assurance that the tool has not caused some unexpected problem.

An advantage with using PSNR is that it is mathematically simple and therefore straightforward to optimize for. As an example, if a tool depends on filter coefficients or other parameters, the reference encoder can search for the parameter value that minimizes Mean Square Error (MSE) and thus optimizes PSNR, and this type of optimization is often straightforward to analyse and implement. The idea is that a real encoder can optimize for a different distortion metric that is psychovisually more relevant but where the parameter search may be a lot more complicated to implement. By choosing PSNR as the distortion metric in our BD-rate calculations, the work can concentrate on creating coding tools instead of spending time developing encoder optimizations for advanced distortion metrics.

For high-dynamic range (HDR) / wide colour gamut (WCG) material and 360° video material, there are additional aspects that influence the usability of BD-rate calculations; these are addressed in clauses 7 and 8, respectively. The JVET common test conditions also specify a separate category for screen content material (i.e., material that has not been captured by a camera). However, in the context of standardization development, the group has not yet seen a need for a special metric and is still using PSNR-based BD-rate for this category.

7 PSNR-based BD-rate calculation

7.1 General

There are several steps in the BD-rate calculation process, where the result in each step is calculated from the result obtained in the previous step:

1. Calculation of PSNR for individual frames
2. Calculation of per-sequence PSNR and bit rate values for each quantization parameter (QP) value. The QP value influences the resulting bit rate. Hence, compressing the sequence several times with different QPs ensures that the final BD-rate measurement will reflect the performance at many different bit rates.
3. Calculation of per-sequence BD-rate values
4. Calculation of an aggregate BD-rate value for all sequences.

These steps are further described in the subclauses below.

7.2 Calculation of PSNR for individual frames

For an individual frame, the mean square error is calculated between the luma channel $decY$ of the decoded output image and the luma channel $origY$ of the original image according to

$$MSE_Y = \frac{1}{W * H} \sum_{y=0}^{H-1} \sum_{x=0}^{W-1} (decY(x, y) - origY(x, y))^2, \quad (1)$$

where $decY(x, y)$ and $origY(x, y)$ are the luma sample values at position (x, y) of the decoded and original images at the same time instance, respectively. W and H are the width and height of the luma component, respectively. The luma PSNR value for the frame is then calculated as

$$PSNR_Y = 10 * \log_{10} \left(\frac{(255 \ll (\text{bitDepth} - 8))^2}{MSE_Y} \right), \quad (2)$$

where $\text{bitDepth} = 10$ for 10-bit inputs and where \ll denotes a bitwise left-shift operation. If $MSE_Y = 0$, i.e., if the decoded image exactly matches the original image, there is some adjustment applied to avoid a division by zero. Different implementations may use a different adjustment method. For example, the HEVC test model (HM) and VVC test model (VTM) software packages set the $PSNR_Y$ value to 999.99 in this case, the AVC test model (JM) and HDRTools software packages impose a minimum MSE of $1 \div (W * H)$, and another approach could be to impose a minimum MSE of $1 \div 12$, since that is the MSE that would theoretically result from rounding large numbers to the nearest multiple of 1 LSB. The use of $255 \ll (\text{bitDepth} - 8)$ instead of $2^{\text{bitDepth}} - 1$ is a small adjustment so that if the same video content is coded using $\text{bitDepth} = 8$ or is coded by shifting it up by two bits and using a 10-bit encoder, and when any error is also just scaled up accordingly, there will be no difference in the resulting fidelity measurement. The difference between the two types of measurement is just a constant offset of 0.0255 dB, so it is normally insignificant.

Three PSNR numbers are ordinarily calculated in this manner; one for luma ($PSNR_Y$), and two for chroma ($PSNR_U$ and $PSNR_V$).

7.3 Calculation of sequence PSNR and bit rate numbers for each QP value

The aggregate PSNR for a test sequence is calculated as the average of the PSNR values for the individual frames:

$$PSNR_Y_{sequence} = \frac{1}{NumFrames} \sum_{k=0}^{NumFrames-1} PSNR_Y_k, \quad (3)$$

where $PSNR_Y_k$ is the $PSNR_Y$ value for frame k calculated according to the subclause 7.2 and $NumFrames$ is the number of frames in the sequence. An alternative to averaging PSNR would be to average the MSE value and then use Equation 2 to calculate the aggregate PSNR for the sequence. That would avoid the issue with dividing by a zero MSE_Y value in Equation 2 when a single decoded frame matches the original perfectly. More generally, it would avoid the case where a single frame with very high fidelity has a large influence on the average. However, that would also mean that a single frame with very poor fidelity could influence the final number considerably, although it may arguably be difficult to notice, especially at high frame rates. It has been the typical practice to average the PSNR scores instead. The bit rate for the sequence is calculated in kilobits per second and is calculated from the number of frames per second (fps), the number of frames in the sequence, and the size of the file in bytes according to:

$$BitRate = \frac{8 * FilesizeInBytes * fps}{NumFrames * 1000}. \quad (4)$$

It should be noted that there is sometimes extra information in the bitstream, such as checksums, that is not necessary for decoding. This information is only used for bitstream validation and is not counted in FilesizeInBytes. Each test sequence is compressed using four different QP values (values 22, 27, 32 and 37 according to the JVET common test conditions). PSNR numbers and bit rate numbers are calculated for each QP.

For chroma, $PSNR_{U_{sequence}}$ and $PSNR_{V_{sequence}}$ are calculated in a similar fashion.

7.4 Calculation of sequence BD-rate number

Subclauses 7.2 and 7.3 have determined the PSNR values and bit rate values for each QP value, both for the anchor and for the tested method. The anchor here refers to the baseline that a tested method is compared against, such as the HEVC reference software HM-16.20, whereas the test is the tested method under investigation, for instance the VVC reference software VTM-7.0. **Error! Reference source not found.** presents one example of how the values may differ between the anchor and test scenarios:

Table 1 – Example bit rates and PSNR values for anchor and test

QP value	Bit rate of anchor (kbps)	PSNR_Y anchor	Bit rate of test (kbps)	PSNR_Y test
22	29419.76	40.19	28020.45	40.38
27	8876.16	39.44	7622.83	39.70
32	4564.60	38.42	3661.62	38.86
37	2551.37	36.90	1979.02	37.54

The values in this table can be plotted as two curves as shown in Figure 1 and Figure 2.

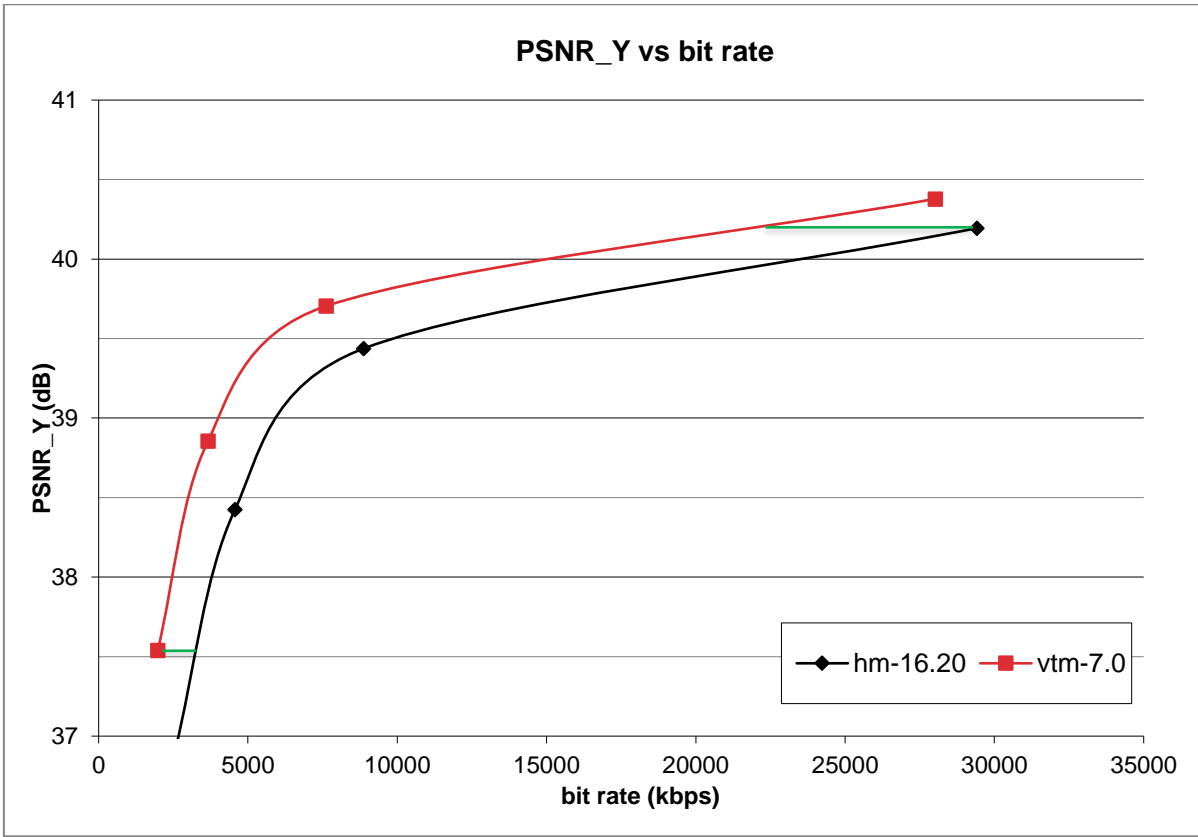


Figure 1 – The luma PSNR plotted as a function of bit rate of HM-16.20 (black) vs VTM-7.0 (red).

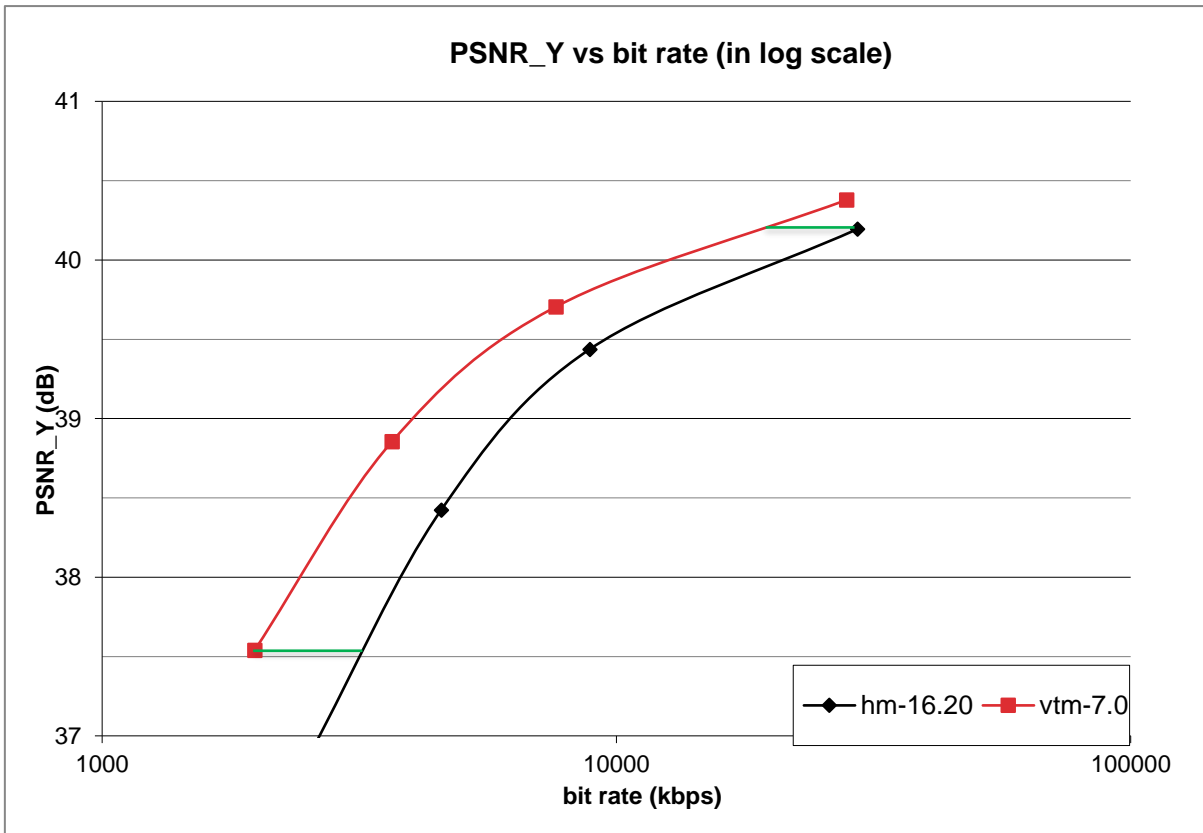


Figure 2 – The luma PSNR plotted as a function of the bit rate of HM-16.20 (black) vs VTM-7.0 (red), where the bit rate axis is in log scale, as suggested in [3].

The idea is now to estimate the area between these two curves to compute an average bit rate savings for equal measured quality, i.e., the Bjøntegaard-delta rate (BD-rate) [1] **Error! Reference source not found.** This is done with the help of a piecewise cubic fitting¹ of the PSNR-Y/rate curves, where the bit rate is measured in the log domain. An integration is then performed to determine the area between the two curves. The area is further divided by the PSNR range to obtain an average rate difference. The details of how this is done are described in references [1][2][3][4][5][8] and [9]. Currently a VBS script `bdrate()` from the Excel file available in [9] is used. A VBS script `bdrateOld()` is also available and computes a BD-rate value using cubic fitting, which was used historically. The two BD-rate values obtained with piecewise cubic and cubic fitting may diverge significantly, indicating a numerical instability. When the two values are close to each other, the result is considered to be more reliable. There are also newer BD-rate functions in [5] that can handle more than 4 points, can also calculate the BD-rate within a specified range, or can also consider extrapolation techniques to deal with cases with minimal overlap between the two curves.

The PSNR ranges of the curves generally do not overlap completely, and the current practice is to avoid measuring the part of the area between curves that are extrapolated² rather than interpolated, which can give unpredictable results if the non-overlapping parts are large. Therefore, the area between the two curves is only measured in the region where there is an overlap, i.e., the area between the two green lines in Figure 1. In the above example, this would mean that BD-rate is calculated only within the area between `minPSNR` and `maxPSNR`, where

`minPSNR = Max(Mini = 0..N(PSNR_anchori), Mini = 0..N(PSNR_testi)) = 37.5394` in this example, and

`maxPSNR = Min(Maxi = 0..N(PSNR_anchori), Maxi = 0..N(PSNR_testi)) = 40.1949` in this example.

This is done by a script provided in [9].

Calculating the value only where there is overlap poses another challenge; if the overlap is only in a very small region, the BD-rate will be calculated using only a small (and possibly atypical) part of the available data. Therefore, it is important that the overlap is substantial for the BD-rate value to be meaningful.

Another issue to be careful of is if the shapes of the curves are very complicated or have unusual characteristics. Especially if they are crossing over each other multiple times, the BD-rate value can be unreliable.

A negative value indicates a gain, i.e., an improvement in coding efficiency. As an example, if the luma BD-rate value is -1.0% , this means that it is possible to compress using the "test" method using 1% fewer bits than using the "anchor" method while maintaining the same luma PSNR.

An alternative approach to calculating BD-rate is to calculate BD-PSNR. Instead of giving an answer like "at equal PSNR, the tested method has 1% lower bit rate", BD-PSNR gives an answer like "at equal bit rate, the tested method has 0.05 dB higher PSNR". However, while 1% lower bit rate at a given quality is a simple concept to understand, it may not be immediately obvious what an increase in PSNR of 0.05 dB means. Hence the common practice is to use BD-rate.

¹ A piecewise-cubic Hermite interpolating polynomial is used.

² It should be noted that no extrapolation occurs when using piecewise cubic fitting. However, extrapolation occurs with cubic fitting.

7.5 Consideration of chroma fidelity

Whereas the PSNR values will be different for luma and chroma, the same bit rate is used in both cases, since the encoding represents the three components together and it is not very feasible to try to separate which bits to assign to which components. Since more of the bits are used to encode luma channels than are used for the chroma channels, this means that the chroma BD-rate values can become difficult to interpret if they deviate too much from the luma BD-rate values. If the BD-rate measures are very different for the luma and chroma components or have opposite signs, the results can be misleading. As an example, if the luma BD-rate value is +0.5% (Y), while the chroma differences are -10.0% (U) and -9.0% (V), it may be difficult to judge which method (test or anchor) is actually better in terms of compression efficiency. A common way around this problem is to carry out a new test where bits are transferred from chroma to luma, for instance by increasing the step-size used for chroma quantization. If the new results are -0.5% (Y), -0.03% (U), -0.02% (V), it is safer to say that the tested method is better. An alternative approach that provides a rough simplified measurement and does not require running a new simulation is to calculate a weighted per-sequence combined PSNR average, for example:

$$PSNR_{YUV_{sequence}} = \frac{1}{8} (6 * PSNR_{Y_{sequence}} + PSNR_{U_{sequence}} + PSNR_{V_{sequence}}) \quad (5)$$

Here the stronger weighting of luma PSNR is to compensate for the fact that most of the bits are used to describe luma information. The per-sequence $PSNR_{YUV_{sequence}}$ values are then used together with the bit rate values to obtain a YUV-BD-rate value for each sequence. These YUV-BD-rate numbers can be helpful especially if there are many methods that should be compared with each other. Another possibility, which is not recommended, is to simply create an average of the BD-rates. In such case a larger weight is typically used for the luma channel's BD-rate. As an example, a BD-rate difference of +0.5% (Y), -10.0% (U), -9.0% (V) would be averaged to $(6 * 0.5\% - 10.0\% - 9.0\%) / 8 = -2.0\%$. This other method may misrepresent gains by a substantial amount when the gain in the chroma channels is substantially larger than the gain in the chroma channel.

7.6 Calculation of aggregate BD-rate value for all sequences

Once the BD-rate values for a set of test sequences have been determined, they are typically combined using an arithmetic average:

$$\text{BD-rate aggregated over several sequences} = \frac{1}{N} \sum_{k=1}^N \text{BD-rate for test sequence } k \quad (6)$$

Typically, the test sequences are divided into classes that are categorized mainly by resolution or by other characteristics, such as whether they contain camera-captured content versus containing text and graphics with motion. One aggregate BD-rate value per class is reported. Hence one number is reported for all the HD sequences, one for WVGA resolution sequences, and so on. Finally, one aggregate BD-rate value may be calculated by averaging across sequences of different classes.

8 BD-rate calculation for HDR material

As discussed in the introduction, a PSNR-based BD-rate measurement is limited in its capability to predict subjective quality improvements. For high-dynamic range (HDR) material, there is a further complication in that there is a very non-linear mapping between the luma codewords that are used as an input to the encoder, and the luminance values that would be output by a display. This is especially true for content that is represented with the Perceptual Quantizer (PQ) transfer function [10][11]. This transfer function gives a much higher importance and allocates a large number of codewords to darker regions, as compared to transfer functions typically used for standard dynamic range content. In our experience, this reduces the correlation between the PSNR-based BD-rate measurement and subjective quality.

For HDR content, it is the experience of the JVET and JCT-VC groups that calculating the PSNR on the luma and chroma codewords, as is done for the standard dynamic range (SDR) case, is still appropriate for HDR content that employs an Hybrid Log-Gamma (HLG) transfer function [11]. However, as mentioned, it is necessary to complement the PSNR metric with a number of additional metrics for HDR content represented with the PQ transfer function [7]. These metrics include:

- PSNRL100: This metric is calculated using the luminance values rather than the luma codewords. It is based on the CIELab colour space representation of the input and output sample values of the codec.
- wPSNR: This is a PSNR-like metric calculated from the codewords, that attempts to compensate for the more significant distribution of luma codewords to the darker regions mentioned above by performing a weighting of the codewords before calculating the PSNR.
- DE100: This is a metric based on the CIELab colour space representation of the input and output sample values of the codec. It is specifically targeted at chrominance fidelity.

For more details on these metrics, see [7].

It is also noted that the JVET and JCT-VC groups do not typically create combined metrics using the wPSNR information from the luma and chroma channels, as was described in clause 7. Instead, the group relies on the PSNRL100 and DE100 metrics. If a combined metric is desired, then attention should be given to the colour gamut used for the content. In particular, the weighting would need to be adjusted when using the ITU-R BT.2100 colour gamut [11] typically employed for HDR content. The JVET and JCT-VC groups have not codified any specific adjustment to date.

9 BD-rate calculation for 360° video

For 360° omnidirectional video applications, the video scene is theoretically positioned on the interior of a distant sphere surrounding the viewer, who is looking out from the centre of the sphere. Since video encoding ordinarily works by compressing rectangular arrays of video samples rather than spheres, a projection must be used to map the scene content on the sphere to the values of the samples of a rectangular array before compression can take place. This also means that measuring the PSNR value of the rectangular video array can be very misleading, since the amount of area on the sphere that corresponds to the sample values at different positions can be very different. As an example, one possible projection is the equirectangular projection, which is similar to a simplistic mapping between the geography depicted on a globe and on a 2D map of the world. If the PSNR value were to be calculated on the rectangular video that was created using such a mapping, it would result in a strong over-emphasis of errors near the poles, since these regions are stretched out considerably. In order to get around this problem, a modified measurement known as the Weighted to Spherically uniform PSNR (WS-PSNR) has been used. This is a weighted form of PSNR that compensates for this stretching. For more information on WS-PSNR, see [6].