ISO/IEC JTC 1/SC 29/WG 1
(ITU-T SG16)

# Coding of Still Pictures

**JBIG**
Joint Bi-level Image
Experts Group

**JPEG**
Joint Photographic
Experts Group

**TITLE:**      Call for Evidence on Learning-based Image Coding Technologies (JPEG AI)

**SOURCE:**      WG1

**PROJECT:**

**STATUS:**      Approved

**REQUESTED ACTION:**      For dissemination

**DISTRIBUTION:**      Public

**Contact:**
ISO/IEC JTC 1/SC 29/WG 1 Convener – Prof. Touradj Ebrahimi
EPFL/STI/IEL/GR-EB, Station 11, CH-1015 Lausanne, Switzerland
Tel: +41 21 693 2606, Fax: +41 21 693 7600,  E-mail: Touradj.Ebrahimi@epfl.ch

# Call for Evidence on
# Learning-based Image Coding Technologies (JPEG AI)

**Summary**

The JPEG Committee has launched the learning-based image coding activity, also referred to as JPEG AI. This activity aims to find evidence for image coding technologies that offer substantially better compression efficiency than available image codecs with models obtained from a large amount of visual data and that can efficiently represent the wide variety of visual content that is available nowadays.

This document is the Call for Evidence (CfE) and has been issued as outcome of the 86th JPEG meeting, Sydney, Australia, 20-24 January 2020. The CfE is designed to be in coordination with the IEEE MMSP 2020 Challenge on Learning-based Image Coding and will use the same content, evaluation methodologies and deadlines.

The deadline for registration is May 30th, 2020. Submissions to the Call for Evidence are due 12th June, 2020 (decoder) and 18th June (code-streams and decoded images).

## 1. Introduction

Image coding algorithms create compact representations of an image by exploiting its spatial redundancy and perceptual irrelevance, thus exploiting the characteristics of the human visual system. Recently, data driven algorithms such as neural networks have attracted a lot of attention and have become a popular area of research and development. This interest is driven by several factors, such as recent advances in processing power (cheap and powerful hardware), the availability of large data sets (big data) and several algorithmic and architectural advances (e.g. convolutional layers).

Nowadays, neural networks are the state-of-the-art for several computer vision tasks, such as those requiring high-level understanding of image semantics, e.g. image classification, object segmentation, saliency detection, but also low-level image processing tasks, such as image denoising, inpainting and super-resolution. These advances have led to an increased interest in applying deep neural networks to image coding, which is the main focus of the JPEG AI ad hoc group within the JPEG standardization committee. The aim of these novel image coding solutions is to design a compact image representation model that has been obtained (learned) from a large amount of visual data and can efficiently represent the wide variety of visual content that is available today. Some of the early learning-based image coding solutions already show encouraging results in terms of rate-distortion performance [1], notably in comparison with conventional image codecs (e.g. JPEG 2000 and HEVC Intra) which code the image information with hand-crafted transforms, entropy coding and quantization schemes.

The CfE is designed to be in coordination with the IEEE MMSP 2020 Challenge on Learning-based Image Coding and will use the same content, evaluation methodologies and deadlines. The intention is to allow academic organizations who do not usually contribute to standardization to have an opportunity to compare their image coding algorithms to those submitted for standardization and likewise to allow actors in standardization to inform academics of their submitted solutions.

## 2. Scope

This Call for Evidence (CfE) on Learning-based Image Coding Technologies solicits technical contributions that demonstrate efficient coding of image content based on a learning-based approach. This means image coding solutions for which learning-based modules (end-to-end trained) play a central role, considering the overall codec architecture. Some examples are non-linear data correlation transformations, probability modeling for entropy coding, hierarchical coding structures, perceptual optimizations and so on. This CfE does not call for conventional coding solutions (e.g. based on HEVC) where neural networks are only used for the optimization of specific modules.

Considering the above context, the main objective of this CfE is to objectively and subjectively evaluate relevant learning-based image coding solutions to demonstrate the potential of this coding approach, especially in terms of compression efficiency. This topic has received many contributions in recent years and

is considered critical for the future of image coding. Naturally, improvements on some aspects (mainly tools) of existing learning-based image codecs are also welcome. The image coding solution should at least support images with the following attributes:

- Image resolution: from 256x256-size images up to 8K images.
- Different types of content, including natural, synthetic, and screen content.
- Bit depth: 8-bit.
- Color space: RGB; YCbCr and ICtCp are optional.
- Input type of the encoder shall match output type of the decoder, this means RGB.
- Internal color space conversion is permitted but should be documented. The same applies for chrominance subsampling.

## 3. Submission requirements

Proponents are asked to submit detailed technical description of the entire image codec, decoding algorithm implementation in software and the decoded test images. Proponents are expected to compress some test material to be provided, with their codec using some target bitrates (see Annex B.3). In all cases, participants are required to submit material to validate the performance of their submission according to the procedure outlined next, notably:

- A detailed description of the coding algorithm, methodologies as well as data used for training, as compression performance alone is not the only evaluation criterion. This description can take the form of an MMSP 2020 submission or a technical report if the proponents do not desire to make a paper submission. Refer to Annex A for more information.
- A decoder implementation in a form that allows stand-alone inference/testing on a standard computer (CPU only) in a reasonable amount of time, preferably in source code form. The PNG still image format should be used for the decoder output. Refer to Annex A for more information.
- Compressed codestreams.
- Corresponding decoded images.

Decoded images should be sRGB with 8 bits per component. Contributors are also expected to provide to JPEG sufficient rights to allow usage of the provided software for the purpose of evaluation. The evaluation process may need to crop and/or clip the provided images to make them suitable for subjective evaluation. The submission requirements of the proposed solution are detailed in Annex A.

## 4. Submission Assessment Criteria

A set of common test conditions (CTC) for learning-based image codecs was designed in the context of the JPEG AI ad hoc group, which defines training and testing datasets, benchmarking codecs, coding conditions (especially target bitrates) and a set of reliable objective quality metrics and subjective assessment procedures [2]. This CTC allows to exhaustively evaluate multiple aspects of learning-based image codecs to fully understand their strengths and weaknesses, notably regarding already available image coding technology.

The CfE evaluation process will be based on the JPEG AI CTC and will focus mostly on compression performance. There are three types of criteria that will be used for the evaluation of the submissions:

- Novelty of the technical contributions with respect to state-of-the-art and relevance (or fitness) to the Call for Evidence objectives. Note that compression performance alone is not the only aspect for the assessment of an image coding technology and other factors will be considered such as bitrate control, number of models required to cover the target bitrates (this means a wide range of qualities) and encoding and decoding complexity.
- Objective evaluation with quality metrics including at least MS-SSIM, VMAF, VIFP, NLPD, FSIM which have been shown to provide high correlation [1] with perceptual assessment for learning-based image codecs.
- Subjective evaluation to be performed with a double stimulus (DSIS) protocol, thus collecting MOS values. These tests will be made in a controlled environment, following well-defined procedures established by ITU standards.

For the subjective testing to be manageable, it is possible that not all submitted proposals will be selected for final subjective performance evaluation. The selection will involve a first assessment stage using a set of objective metrics which should allow to find the top performing submissions. After the identification of the highest RD performance submissions (by measuring distortions objectively) and a review of the technical contributions and relevance of each one considering the Call for Evidence objective, some proposals will be selected to be evaluated with an appropriate subjective evaluation procedure (2nd stage).

From the subjective assessment tests, the top performing solutions will be drawn, not only considering the subjective performance (including confidence intervals, target bitrates for which have been superior) but also the novelty, complexity and relevance regarding the Call for Evidence objective and complexity.

In addition to the selected submitted image codecs, the quality evaluation will also include well-known standardized image coding solutions such as JPEG, JPEG 2000 and HEVC.

## 5. Timeline

The intended timeline for the evaluation of the proposals is the following:

| | |
|---|---|
| **March 2nd, 2020** | Website created and online. Challenge and Call for Evidence announcement. |
| March 9th, 2020 | Release of the training and validation parts of the dataset. |
| **May 30th, 2020** | Proposal registration. |
| June 12th, 2020 | Submission of decoder implementation with some fixed model. No (re)training is allowed after this date. |
| June 14th, 2020 | Release of the test dataset for proponents to code. |
| **June 18th, 2020** | Submission of code-streams and decoded images for the test dataset. |
| June 21st, 2020 | MMSP paper and/or JPEG technical document submission. |
| July 15th, 2020: | MMSP paper notification. |
| **July 30th, 2020** | MMSP camera ready paper. |
| September 1st, 2020 | 1$^{st}$ Stage: Objective evaluation of all the proposals (1st stage); results will be released online, including which proposals will be subjectively evaluated <br><br> 2$^{nd}$ Stage:  Subjective evaluation. |
| September 15th, 2020 | End of subjective evaluation; results will be released online. |
| MMSP 2020, 21st-23rd September, Tampere, Finland | Challenge session at MMSP 2020 with presentations and announcement of awardees. |
| **89th JPEG Meeting** | Presentation and discussion of the proposals at JPEG meeting. Attendance is recommended for proponents. |

All information should be submitted electronically in a place to be provided in the website created for the Call for Evidence (and MMSP) Challenge.

## 6. Datasets

The JPEG AI database was constructed to (i) evaluate the performance of state-of-the-art learning-based image coding solutions and (ii) to be used for training, validation and testing of learning-based image coding solutions. This dataset will be made available to all participants of the CfE. The JPEG AI dataset will be organized according to:

- **Training dataset:** The training dataset aims to provide a set of images to create a model suitable for a learning-based image codec solution. However, the proponents may also use a different training

dataset provided that it is fully identified in the proposal descriptions and, ideally, made available for future developments.

- **Validation dataset:** The validation dataset aims to provide a set of images to be used during the training to validate the convergence of the training algorithm employed by some learning-based image codec solution.

- **Test dataset (hidden):** The test dataset cannot be used neither for training or for validation and will be used to evaluate the final performance of learning-based image coding solutions. Test images are kept hidden until some appropriate stage, to avoid being used for training. The test dataset for the evaluation of the Call for Evidence (and Challenge) submissions will be drawn from a sizeable repository that is maintained by JPEG.

The diversity of the images contained in the JPEG AI dataset is high, namely in terms of their characteristics, such as content and spatial resolution. These datasets have the following characteristics:

- Format – PNG images (RGB color components, non-interlaced);
- Spatial resolution – from 256×256 to 8K (8 bit);
- Training/validation/test dataset: 5264/350/25 images.

The number of images allows for an efficient training/validation and is typically larger than the number of images used in previously available works. The number of test images provides a well-balanced set of diverse images that can be used for the evaluation of learning-based image coding solutions. The training and validation dataset will be available at sftp://jpeg-cfe@amalia.img.lx.it.pt (password to be given by request) by 9th March, 2020.

## 7. Call for Evidence Details

This CfE invites proponents to submit technology contributions that fulfill the scope and objectives according to the timeline presented above. Proponents are encouraged to attend the 90th JPEG meeting and present their findings.

### 7.1 Submission requirements

A submission shall consist of the elements specified in Annex A. All the elements to be submitted, excluding the decoded images, should be uploaded to the WG1 document registry. For the decoded images, instructions will be provided after the registration. Those proponents without access to the registry should contact the WG1 members listed in Section 9.

### 7.2 IPR conditions (ISO/IEC Directives)

Proponents are advised that this call is being made in the framework and subject to the common patent policy of ITU-T/ITU-R/ISO/IEC and other established policies of these standardization organizations. The contact

persons named in Section 9 can assist potential submitters in identifying the relevant policy information.

### 7.3 Royalty-free goal

The royalty-free patent licensing commitments made by contributors to previous standards, e.g. JPEG 2000 Part 1, have arguably been instrumental to their success. JPEG expects that similar commitments would be helpful for the adoption of a future JPEG AI image coding standard.

## 8.  JPEG AI e-mail reflector information

E-mail reflector: jpeg-ai@jpeglists.org

To subscribe to the reflector, please visit http://jpeg-ai-list.jpeg.org or in case of problems contact lists@jpeg.org

## 9.  Contacts

Touradj Ebrahimi (JPEG Convener)
Email: Touradj.Ebrahimi@epfl.ch

Fernando Pereira (JPEG Requirements Chair)
Email: fp@lx.it.pt

João Ascenso (AhG JPEG AI Chair)
Email: joao.ascenso@lx.it.pt

# ANNEX A – SUBMISSION REQUIREMENTS

The process to evaluate proposals will be done following the timeline defined in Section 3. In addition to the technical description and other elements (code or binaries, decoded images, etc) to be submitted, proponents are encouraged to contribute to the standardization process, namely by participating in the JPEG AI ad-hoc group.

## A.1. Proposal description

Each proposal must include a technical description of the entire image codec, in the form of a 2-column paper (using the MMSP paper template and submission rules), namely:

- Key features of the proposal, including the target quality range.
- High-level description of the proposal.
- Encoder/decoder architecture, training procedure including loss function used, strategies on how to deal with the non-differentiable quantization and bitrate allocation.
- Model size if applicable, which corresponds to the number of weights (and precision of each weight) in the encoder and decoder.
- RD performance for at least four test images (to be specified) using MS-SSIM and VMAF objective metrics.
- Running time (encoder and decoder) for some CPU and GPU platform. Include all the details of the platform (CPU model, clock rate and memory, GPU model and brand) but also the deep learning framework (e.g. Tensorflow or Pytorch). The recommended platform for GPU is NVIDIA 2080 Ti, but you may use other one.

## A.2. Additional elements

The following additional elements must be submitted by all proposals:
- Standalone executable package: docker file with all the libraries and tools to run the decoder with the submitted code-streams and preferably decoder in source code form. All the information to run the decoder shall be provided. If binaries are used, they should correspond to statically linked Linux executables with all required libraries and system dependencies
- Code-streams corresponding to the encoded test images to be used for decoding.
- Decoded test images for objective and subjective evaluation. All test images will be made available to proponents on the website created for the Call for Evidence (and Challenge).
- Training dataset (if the JPEG AI dataset was not used). This is optional.

## ANNEX B – EVALUATION PROCEDURES

### B.1. Anchors

Proposals will be compared against the following anchors:

- JPEG (ISO/IEC 10918-1 | ITU-T Rec. T.81)
- JPEG 2000 (ISO/IEC 15444-1 | ITU-T Rec. T.800)
- HEVC Intra (ISO 23008-2 | ITU-T Rec. H.265)

Information on available software and configurations to be used for these anchors is given in Annex C.

### B.2. Evaluation procedures

Objective and subjective quality evaluation of the proposals will each be done by at least two independent labs, following well-established procedures and based on the decoded test images provided by each proponent. The submitted code (or binaries) for the decoder, codestreams and decoded images will be used for verification purposes. In Figure 1, the coding pipeline for learning-based image coding solutions, which is rather straightforward, is presented. Proponents may perform encoding with any color space representation, but the input of the encoder and the output of the decoder must be in the PNG (RGB color space) format. Objective image quality will be measured with luminance and color-based metrics and the RGB decoded images will be used for quality evaluation.
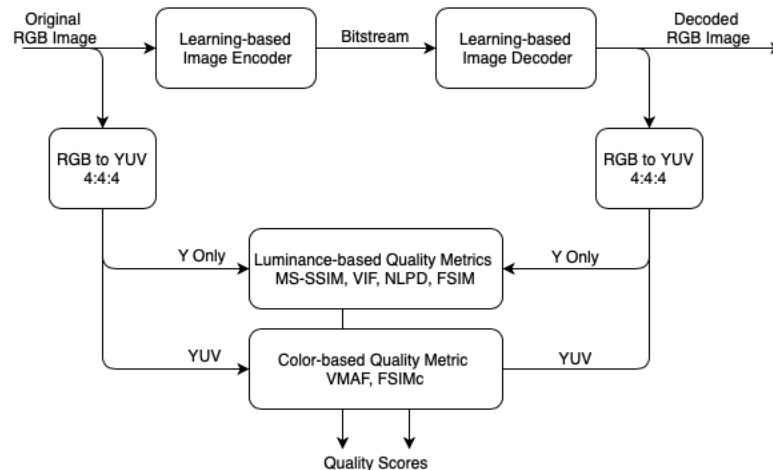


*Figure 1 – Encoding-decoding pipeline for learning-based image coding solutions.*

### B.3. Target rates

Target bitrates for the objective evaluations include 0.06, 0.12, 0.25, 0.50, 0.75, 1.00, 1.50, and 2.00 bpp. The maximum bitrate deviation from the target bitrate should not exceed 15%. The proponents must declare for every test image which target bitrate their decoder and models can reach, and in case of deviation of the target bitrate, the proposed RD point may not be considered for evaluation. The target bitrates for the

subjective evaluations will be a subset of the target bitrates for the objective evaluations and will depend on the complexity of the test images.

The bitrates specified should account for the total number of bits necessary for generating the encoded file (or files) out of which the decoder can reconstruct a lossy version of the entire image. The main rate metric is the number of bits per pixel (bpp) defined as:

$$BPP = \frac{N\_TOT\_BITS}{N\_TOT\_PIXELS}$$

where N_TOT_BITS is the number of bits for the compressed representation of the image and N_TOT_PIXELS is the number of pixels in the image.

## B.4. Objective quality testing

Objective quality testing shall be done by computing several quality metrics, including MS-SSIM, VMAF, VIFP, NLPD, FSIM, between compressed and original image sequences, at the target bitrates mentioned in the precious Section. Refer to Annex C for more information about the image quality metrics.

## B.5. Subjective quality testing

To evaluate the selected coding solutions, a subjective quality assessment methodology may be used. This type of assessment is especially critical since the type of artifacts that learning-based image compression solutions produce may be significantly different from those in standard image codecs. Subjective quality evaluation of the compressed images will be performed on the test dataset.

The Double Stimulus Impairment Scale (DSIS) methodology will be used, where subjects watch side by side the original image and the impaired decoded image which is scored in a 1-5 scale associated to five impairment scale, notably very annoying, annoying, slightly annoying, perceptible but not annoying and imperceptible. The side-by-side images are centered on the display. The reference is shown on the left or right and can vary randomly. Subjects will be informed that one of the images is the reference but will not receive any indication whether the reference image was on the left or on the right.

The subjective test methodology will follow BT500.13 [3] and a randomized presentation order, as described in ITU-T P.910 [4] will be used; the same content is never displayed consecutively. There is no presentation or voting time limit. A training session should be organized before the experiment to familiarize participants with artefacts and distortions in the test images. At least, three training images will be used before actual scoring.

As anchors, JPEG, JPEG 2000 and HEVC will be used. The list of anchors may be reduced if the number of proposals is too high. The images used for subjective evaluation are a subset of the test dataset images and its number will be selected depending on the number of proposals that will be subjectively evaluated.

# ANNEX C – ANCHOR CONFIGURATION

The configurations detailed below are relevant for the definition of anchors.

## C.1. JPEG (ISO/IEC 10918-1 | ITU-T Rec. T.81))

JPEG does not specify a rate allocation mechanism allowing to target a specific bitrate. Hence, an external rate control loop is required to achieve the targeted bitrate. The following conditions apply:

- Software to be used: JPEG XT reference software, v1.53
  - Available at http://jpeg.org/jpegxt/software.html.
  - License: GPLv3
- Command-line examples (to use within the rate-control loop): jpeg -q [QUALITY_PARAMETER] [INPUTFILE] [OUTPUTFILE]

## C.2. JPEG 2000 (ISO/IEC 15444-1 | ITU-T Rec. T.800)

The JPEG 2000 anchor generation should support two configurations: 1) PSNR optimized; and 2) Visually optimized. A target rate can be specified using the –rate [bpp] parameter. The following conditions apply:

- Software to be used: Kakadu, v7.10.2
  - Available at http://www.kakadusoftware.com.
  - License: demo binaries freely available for non-commercial use
- Command-line examples:
  - **MSE weighted:** kdu_compress -i [INPUTFILE] -o [OUTPUTFILE] -rate [BPP] Qstep=0.001 -tolerance 0 -full -precise
  - **Visually weighted:** kdu_compress -i [INPUTFILE] -o [OUTPUTFILE] -rate [BPP] Qstep=0.001 -tolerance 0 -full -precise -no_weights
  - **Decoding:** kdu_expand -i [INPUTFILE .mj2] -o [OUTPUTFILE .yuv] -precise

## C.3. HEVC (ISO 23008-2:2018 | ITU-T Rec. H.265 (v5))

For HEVC, an external rate control loop is required to achieve targeted bitrate. The HEVC RD performance for the target bitrates are obtained with the following conditions:

- Available software: HEVC Test Model HM-16.20+SCM-8.8
  - Available at https://hevc.hhi.fraunhofer.de/
  - License: BSD
- FFMPEG will be used to convert the PNG (RGB) to YUV files following the BT709 primaries.
- Configuration files to be used are available in [5].

## ANNEX D – OBJECTIVE QUALITY METRICS

This section defines the objective image quality metrics that will be used for the assessment of learning-based image coding solutions.

### 9.1.1   MS-SSIM Definition and Computation

Multi-Scale Structural SIMilarity (MS-SSIM) [6] is one of the most well-known image quality evaluation algorithms and computes relative quality scores between the reference and distorted images by comparing details across resolutions, providing high performance for learning-based image codecs [1]. The MS-SSIM [6] is more flexible than single-scale methods such as SSIM by including variations of image resolution and viewing conditions. Also, the MS-SSIM metric introduces an image synthesis-based approach to calibrate the parameters that weight the relative importance between different scales. A high score expresses better image quality.

The source code of this metric can be downloaded at this link:
https://ece.uwaterloo.ca/~z70wang/research/iwssim/.

### 9.1.2   VMAF Definition and Computation

The Video Multimethod Assessment Fusion (VMAF) metric [7] developed by Netflix is focused on artifacts created by compression and rescaling and estimates the quality score by computing scores from several quality assessment algorithms and fusing them with a support vector machine (SVM). Even if this metric is specific for videos, it can also be used to evaluate the quality of single images and has been proved that performs reasonably well for learning-based image codecs [1]. Since the metric takes as input raw images in the YUV color space format, the PNG (RGB color space) images are converted to the YUV 4:4:4 format using FFMPEG (BT.709 primaries). A higher score of this metric indicates better image quality.

The source code of this metric can be downloaded at this link:
https://github.com/Netflix/vmaf

### 9.1.3   VIF Definition and Computation

The Visual Information Fidelity (VIF) [8] measures the loss of human perceived information in some degradation process, e.g. image compression. VIF exploits the natural scene statistics to evaluate information fidelity and is related to the Shannon mutual information between the degraded and original pristine image. The VIF metric operates in the wavelet domain and many experiments found that the metric values agree well with the human response, which also occurs for learning-based image codecs. A high score expresses better image quality.

The source code of this metric can be downloaded at this link:
https://live.ece.utexas.edu/research/Quality/VIF.htm


### 9.1.4    NLP Definition and Computation

The Normalized Laplacian Pyramid (NLPD) is an image quality metric [9] based on two different aspects associated with the human visual system: local luminance subtraction and local contrast gain control. NLP exploits a Laplacian pyramid decomposition and a local normalization factor. The metric value is computed in the normalized Laplacian domain, this means that the quality of the distorted image relative to its reference is the root mean squared error in some weight-normalized Laplacian domain. A lower score express better image quality.

The source code of this metric can be downloaded at this link:
http://www.cns.nyu.edu/~lcv/NLPyr/


### 9.1.5    FSIM Definition and Computation

The feature similarity (FSIM) metric [10] is based on the computation of two low level features that play complementary roles in the characterization of the image quality and reflects different aspects of the human visual system: 1) the phase congruency (PC), which is a dimensionless feature that accounts for the importance of the local structure and the image gradient magnitude (GM) feature to account for contrast information. Both color and luminance versions of the FSIM metric will be used. A high metric value express better image quality.

The source code of this metric can be downloaded at this link:
https://www4.comp.polyu.edu.hk/~cslzhang/IQA/FSIM/FSIM.htm

## ANNEX E – REFERENCES

[1] J. Ascenso, P. Akayzi, M. Testolina, A. Boev, E. Alshina "Performance Evaluation of Learning based Image Coding Solutions and Quality Metrics", ISO/IEC JTC 1/SC29/WG1 N85013, 85th JPEG Meeting, San Jose, USA, November 2019. Available at:
https://jpeg.org/items/20191203_jpeg_ai_performance_evaluation.html.

[2] J. Ascenso, P. Akayzi "JPEG AI Image Coding Common Test Conditions", ISO /IEC JTC 1/SC 29/WG 1 N84035, 84th Meeting, Brussels, Belgium, 13-19 July 2019.

[3] ITU-R Recommendation BT.500-13, "Methodology for the subjective assessment of the quality of television pictures," International Telecommunications Union, Geneva, Switzerland, 2012.

[4] ITU-T Recommendation P. 910, "Subjective video quality assessment methods for multimedia applications," International Telecommunication Union, Geneva, 2008.

[5] E. Upenik, J. Wassenberg, "JPEG XL Experiment Reproduction", ISO/IEC JTC 1/SC29/WG1 M86083, 86[th] JPEG Meeting, Sydney, Australia, January 2020.

[6] Z. Wang, E. P. Simoncelli and A. C. Bovik, "Multi-scale Structural Similarity for Image Quality Assessment", 37th IEEE Asilomar Conference on Signals, Systems and Computers, Pacific Grove, CA, USA, November 2003.

[7] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy and M. Manohara, "Toward A Practical Perceptual Video Quality Metric", [Online], Available at: https://netflixtechblog.com/toward-a-practical-perceptual-video-quality-metric-653f208b9652

[8] H.R. Sheikh and A. C. Bovik, "Image Information and Visual Quality," IEEE International Conference on Acoustics, Speech, and Signal Processing, Montreal, Canada, August 2004.

[9] L. Zhang, L. Zhang, X. Mou, D. Zhang, "FSIM: a Feature Similarity Index for Image Quality Assessment," IEEE Transactions on Image Processing, vol. 20, no. 8, pp. 2378-2386, August 2011.

[10] V. Laparra, J. Balle´, A. Berardino, and E. P. Simoncelli, "Perceptual Image Quality Assessment using a Normalized Laplacian Pyramid", S&T Symposium on Electronic Imaging: Conf. on Human Vision and Electronic Imaging, San Francisco, CA, USA, February 2016.