

ISO/IEC JTC1/SC29/WG1  
(ITU-T SG21)

## **Coding of Still Pictures**

### **JBIG**

*Joint Bi-level Image  
Experts Group*

### **JPEG**

*Joint Photographic  
Experts Group*

**TITLE:** JPEG AI Use Cases and Requirements v2.0

**SOURCE:** Requirements SG

**PROJECT:** ISO/IEC 6048 (JPEG AI)

**STATUS:** Final

**REQUESTED  
ACTION:** Distribute

**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](mailto:Touradj.Ebrahimi@epfl.ch)

# 1 Introduction

The scope of the JPEG AI is the creation of a learning-based image coding standard offering a **single-stream, compact** compressed domain representation, targeting both **human visualization**, with significant compression efficiency improvement over image coding standards in common use at equivalent subjective quality, and effective performance for **image processing and computer vision tasks**, with the goal of supporting a **royalty-free baseline**.

## 2 JPEG AI Framework

Learning-based image coding solutions have already shown that they can achieve substantially better compression efficiency than existing conventional solutions, namely by exploiting advanced machine learning tools, such as deep neural networks [1]. In particular, it has been shown that, when compared to JPEG, JPEG 2000 and HEVC Intra, learning-based coding solutions can provide better perceptual quality, for some target bitrates, both in terms of appropriate perceptual objective quality metrics and subjective assessment scores [2]. Besides their high compression efficiency, learning-based image coding solutions may be adapted with little extra effort to image processing and computer vision tasks without the need for full decoding, i.e., without performing image reconstruction. This contrasts with classical image codecs, which when used in image processing and computer vision pipelines, need to perform full decoding of the compressed bitstream to obtain a pixel-based representation.

Figure 1 shows the high-level JPEG AI framework, highlighting the three pipelines. The input to the learning-based image coding framework is a digital image and the output bitstream may be processed for human visualization by performing entropy decoding and standard reconstruction, thus producing a standard decoded image. As shown in Figure 1, the standard reconstruction may be skipped since the latent representation produced by the encoder contains the necessary information not only for decoding but also to perform image processing and computer vision tasks at the decoder side (after entropy decoding). These tasks are carried out on the latent representation, directly extracted from the original image and not from the (lossy) decoded image. This intrinsically feature-rich latent representation can be used in two main ways: 1) to perform an image processing task, such as targeting the enhancement or modification of the image, where a processed image is produced, for example with increased resolution, contrast, etc.; and 2) to perform a computer vision task where high-level semantic information is extracted, e.g., to generate classes, labels, regions, etc.

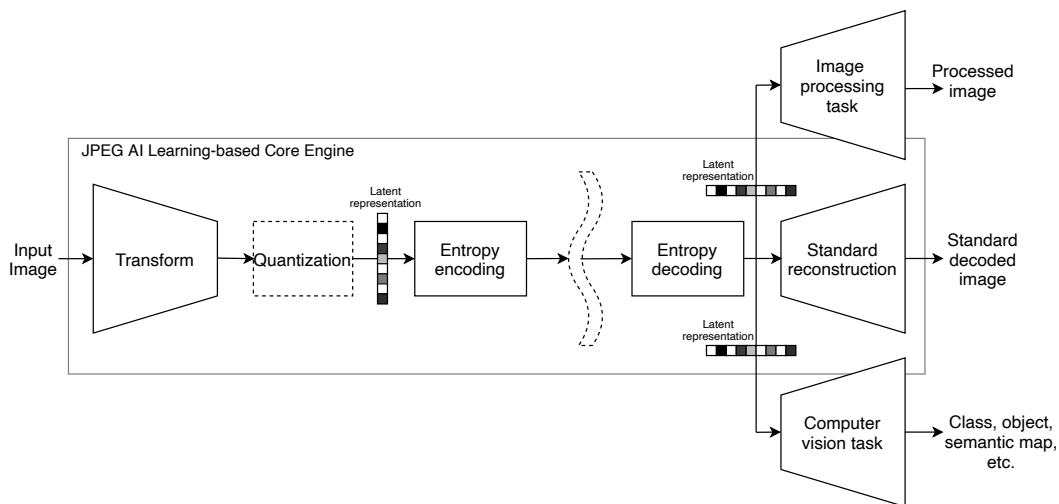


Fig. 1: JPEG AI learning-based image coding framework.

### 3 Key Tasks

Following the JPEG AI scope, the coded bitstream will have a triple-purpose, thus offering compelling advantages for applications where an image processing task aims to enhance or modify the image or where semantic (or higher-level) information needs to be extracted from large amounts of visual data. This may have a significant impact on image processing and computer vision tasks, which may be performed with lower complexity by using as input the compressed domain representation (and thus artifact free) instead the original or decoded images. Some examples of relevant image processing tasks are:

- Super-resolution
- Denoising
- Low-light enhancement
- Color correction
- Exposure compensation
- Inpainting

Relevant examples of computer vision tasks are:

- Image retrieval and classification
- Object detection, recognition and identification
- Semantic segmentation
- Event detection and action recognition
- Face detection and recognition

### 4 Use Cases

This section presents the use cases targeted for the JPEG standard on learning-based image coding. These are used to motivate the requirements that are defined in Section 4.

#### 4.1 Cloud storage

There has been an increasing number of images being stored in the cloud, due to the emergence of several online storage services. Several companies such as Tencent, Microsoft, Facebook and Google, often have thousands of billions of photos stored, which require considerable amount of resources, notably storage space,

bandwidth or energy. Therefore, creating a highly efficient image coding solution for cloud storage is rather important to minimize costs, and even marginal savings in bitrate for some target quality may have a significant impact. The use of learning-based image compression may allow to reduce storage space, thus leading to high quality images at a fraction of the cost. In addition, high compression efficiency allows lower bandwidth costs, which translates into easy transmission and sharing of massive amounts of images.

Key **features** for this application are:

- a. Lossy compression with bit-rate control
- b. High coding efficiency
- c. Perceptual optimization
- d. Efficient compressed domain representation for image processing, e.g. color correction
- e. Privacy preservation

## 4.2 Visual surveillance

Visual surveillance systems are widely deployed to perform video monitoring with several objectives, such as anomaly detection, detection of suspicious activity, provision of forensic evidence and intelligent control. Often, a very large number of cameras generate huge amounts of visual data that needs to be processed, compressed, analyzed and stored. Intelligent surveillance systems are often used to record relevant events not just as video but also as very high-resolution images. In some cases, non-visible light images (and even range maps) are also acquired. Considering the amount of data, the challenge is not only acquiring and compressing visual data but also understanding what is relevant and what can be ignored in an automatic way. Thus, image processing or computer vision tasks are often employed to allow efficient navigation and abnormal activity detection. Examples of such tasks are image search, object detection, crowd behavior analysis and recognition of faces and events.

Key **features** for this application are:

- a. Lossy compression with bit-rate control
- b. High coding efficiency
- c. Fast encoding and decoding
- d. Efficient compressed domain representation for computer vision, e.g. object detection
- e. Efficient compressed domain representation for image processing, e.g. super-resolution
- f. Spatial random access, especially for very high resolution cameras
- g. Offline image enhancement, e.g. for super-resolution
- h. Privacy preservation

## 4.3 Autonomous vehicles and devices

Self-driving cars, drones and other autonomous devices generate a vast amount of visual data that must be analyzed and sometimes stored. Moreover, images collected from autonomous vehicles and devices may need to be processed offline and thus efficiently transmitted and/or stored. For example, drones carry cameras that are programmed to capture several Gigabytes of high-resolution aerial imagery which can be difficult to transmit over resource-constrained connections. Moreover, autonomous driving technology and other automated assistance systems may use several cameras for real-time analysis and decision, but the storage and transmission of key events allows other useful applications, such as traffic monitoring, accident investigation,

etc. This scenario often involves several computer vision tasks, such as object detection, semantic segmentation and event recognition.

Key **features** for this application are:

- a. Lossy compression with bit-rate control
- b. High coding efficiency
- c. Fast encoding and decoding
- d. Efficient compressed domain representation for computer vision, e.g. semantic segmentation
- e. Efficient compressed domain representation for image processing, e.g. low-light enhancement
- f. Lossy to lossless encoding
- g. Privacy preserving

#### 4.4 Image collection storage and management

With the wide deployment of smartphones and other consumer devices, every person has a digital camera which is used to acquire and store images of relevant events in photo albums. This collection of images is often backed up on online web storage to avoid their loss in the event of failure or loss of the smartphone or digital camera. Moreover, since these images usually have very high resolution, they require a significant amount of storage space and their storage has to be organized in a convenient way, to facilitate their search and consumption. In this scenario, image classification, object detection and action recognition can be applied to facilitate the management and organization of images.

Key **features** for this application are:

- a. Lossy compression with bit-rate control
- b. High coding efficiency
- c. Efficient compressed domain representation for computer vision, e.g. image retrieval and classification
- d. Fast encoding and decoding
- e. Privacy preserving

#### 4.5 Live monitoring of visual data

Live streaming of visual data has significantly increased, from professional services such as online lectures, videoconferences and webcasts but also other entertainment services, such as video game live streaming and, short-form personal videos (see snack culture). Often, such visual data has to be analyzed in order to detect inappropriate content (as it is often done in social media networks) that may violate policies but also to provide additional information such as labeling of faces, emotions, gestures and so on. Also, computer vision tasks could be applied to live images/videos to perform intelligent review, rating and distribution of this type of content.

Key **features** for this application are:

- a. Lossy compression with bit-rate control is needed
- b. High coding efficiency
- c. Efficient compressed domain representation for computer vision, e.g. face detection and recognition
- d. Fast encoding and decoding
- e. Privacy preserving

## 4.6 Media distribution

Billions of user-generated images are captured and transmitted over the internet daily. These images are often uploaded and transcoded into multiple quality versions and formats, being stored on worldwide servers for distributions. In such scenario, efficient image compression solutions allow to lower the storage and transmission cost and are especially relevant to users with low-bandwidth wireless connections. Progressive decoding may also be desirable, which allows for useful previews while the image is still being received. This may take the form of lower-resolution versions of images which are sufficient to display in displays with lower resolution, without requiring the resources needed for the entire high-resolution version.

Key **features** for this application are:

- a. Lossy compression with bit-rate control is needed
- b. High coding efficiency
- c. Fast decoding
- d. Perceptual optimization

## 4.7 360° photo sharing

Nowadays many applications support 360 photo sharing, including Facebook and Google Photos. When viewing a 360° image, the user would typically only see the current field of view viewport (portion) of the spheric scene. Moreover, the user may navigate the visual scene by moving his or her head to view other viewports of the scene. Typically, the spherical representation is projected into a 2D image, which is then encoded by an image codec. When the user's viewport changes, a new viewport is requested which largely overlaps with the previous viewport. To enable this functionality in an efficient way regarding both network bandwidth and decoding, the image coded representation needs to be organized into independently decoded regions enabling the transmission of parts of the image and not the full image.

Key **features** for this application are:

- a. Spatial random access
- b. High compression efficiency
- c. Region of interest decoding
- d. Low complexity decoding
- e. Low latency

# 5 Requirements

This section presents the requirements that should be met by the standard so that it can be employed for the above-described use-cases. Requirements are split between “core requirements” which are essential and “desirable requirements” which are not mandatory, but actually might enlarge the possible application scenarios and will be decided depending on their cost.

## 5.1 Uncompressed image attributes

The image coding technology to be standardized *should at least* support images with the following attributes:

- Image resolution: from thumbnail-size images up to 8K, as minimum.
- Bit depth: 8-bit and 10-bit.

- RGB color space (three channels) and monochrome (one channel).
- Different types of content, including natural (photographs, aerial/satellite, document scans and synthetic (illustrations/UI elements/comics).

## 5.2 Compressed bitstream requirements

The standard shall cover at least the core requirements and is encouraged to cover desirable requirements as well.

### Core requirements

- Effective compressed domain image processing and computer vision tasks.
- Significant compression efficiency improvement over coding standards in common use at equivalent subjective quality.
- Reconstructed images with both high subjective quality and high fidelity as measured by full reference objective quality metrics and double stimulus subjective assessment protocols.
- Decoders in different platforms (e.g. CPU vs GPU) shall reconstruct images that are the same (bit exact) for applications which require bit exact reconstructions.
- Decoders in different platforms may have a mismatch in the decoded images which shall be non-perceivable, for applications which do not require bit exact reconstructions.
- Hardware platform agnostic, encoder and decoder should be implementable in a wide range of hardware platforms.
- Hardware/software implementation-friendly encoding and decoding (in terms of parallelization, memory, complexity, and power consumption).
- Support for 8- and 10-bit depth.
- Support for efficient coding of images with text and graphics.
- Support for progressive decoding.

### Desirable requirements

- Support for higher bit depth (e.g., 12 to 16-bit integer and floating-point HDR) images.
- Support for region of interest-based coding.
- Support for progressive decoding up to lossless.
- Support for lossless alpha channel/transparency coding.
- Support for animated image sequences.
- Support for wide color gamut coding.
- Support for different color representations.
- Support for very low file size image coding (e.g. 64×64 pixel images).
- Support for a low-complexity profile - low encode/decode time even on resource-constrained hardware (e.g., mobile devices).
- Minimal generation loss when lossy compression is applied multiple times.
- Support for spatial random access

## 6 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 learning-based image coding standard.