

# Deep Learning-Based Video Coding: A Review and A Case Study

Dong Liu, Yue Li, Jianping Lin, Houqiang Li, Feng Wu

## Abstract

The past decade has witnessed great success of deep learning technology in many disciplines, especially in computer vision and image processing. However, deep learning-based video coding remains in its infancy. This paper reviews the representative works about using deep learning for image/video coding, which has been an actively developing research area since the year of 2015. We divide the related works into two categories: **new coding schemes** that are built primarily upon deep networks (deep schemes), and **deep network-based coding tools (deep tools)** that shall be used within traditional coding schemes or together with traditional coding tools. For deep schemes, **pixel probability modeling and auto-encoder** are the two approaches, **that can be viewed as predictive coding scheme and transform coding scheme**, respectively. For deep tools, there have been several proposed techniques using deep learning to perform **intra-picture prediction, inter-picture prediction, cross-channel prediction, probability distribution prediction, transform, post- or in-loop filtering, down- and up-sampling**, as well as encoding optimizations. According to the newest reports, **deep schemes** have achieved comparable or even higher compression efficiency than **the state-of-the-art** traditional schemes, such as **High Efficiency Video Coding (HEVC)** based scheme, for *image* coding; deep tools have demonstrated the compression capability **beyond HEVC** for *video* coding. **However, deep schemes have not yet reached the current height of HEVC for video coding**, and deep tools remain largely unexplored at many aspects including **the tradeoff between compression efficiency and encoding/decoding complexity, the optimization for perceptual naturalness or semantic quality, the speciality and universality, the federated design of multiple deep tools**, and so on. In the hope of advocating the research of deep learning-based video coding, we present a case study of our developed prototype video codec, namely Deep Learning Video Coding (DLVC). DLVC features two deep tools that are both based on convolutional neural network (CNN), namely CNN-based in-loop filter (CNN-ILF) and CNN-based block adaptive resolution coding (CNN-BARC). Both tools help improve the compression efficiency by a significant margin. **With the two deep tools as well as other non-deep coding tools, DLVC is able to achieve on average 39.6% and 33.0% bits saving than HEVC, under random-access and low-delay configurations, respectively.** The source code of DLVC has been released for future researches.

## Index Terms

Deep learning, image coding, prediction, transform, video coding.

# Deep Learning-Based Video Coding:

## A Review and A Case Study

### I. INTRODUCTION

#### A. Image/Video Coding

Image/video coding usually refers to the computing technology that compresses image/video into binary code (i.e. bits) so as to facilitate storage and transmission. The compression may or may not ensure perfect reconstruction of image/video from the bits, which is termed lossless and lossy coding respectively. For natural image/video, the compression efficiency of lossless coding is usually below requirement, so most of efforts are devoted to lossy coding. Lossy image/video coding solutions are evaluated at two aspects: first is the compression efficiency, commonly measured by the number of bits (coding rate), the less the better; second is the incurred loss, commonly measured by the quality of the reconstructed image/video compared with the original image/video, the higher the better.

Image/video coding is a fundamental and enabling technology for computer image processing, computer vision, and visual communication. The research and development of image/video coding can be dated back to as early as 1960s, much earlier than the appearance of modern imaging, image processing, and visual communication systems. As an example, Picture Coding Symposium, a prestigious international forum devoted specifically to advancements in image/video coding, started in the year of 1969. Since then, numerous efforts from both academia and industry have been devoted to this field.

Due to the requirement of interoperability, a series of standards regarding image/video coding have been crafted in the past three decades. In international standardization organizations, ISO/IEC has two experts group namely Joint Photographic Experts Group (JPEG) and Moving Picture Experts Group (MPEG) for standardization of image/video coding technology, while ITU-T has its own Video Coding Experts Group (VCEG). These organizations have published several famous, widely adopted standards, such as JPEG [121], JPEG 2000 [103], H.262 (MPEG-2 Part 2) [115], H.264 (MPEG-4 Part 10 or AVC) [128], H.265 (MPEG-H Part 2 or HEVC) [108], and so on. At present, H.265/HEVC, which was formally published in 2013, represents the state-of-the-art image/video coding technology.

Along with the progress of video technology, especially the popularization of ultra-high definition (UHD) video, there is an urgent requirement to further increase compression efficiency so as to accommodate UHD video in limited storage and limited transmission bandwidth. Thus, after HEVC, MPEG and VCEG form the Joint Video Experts Team (JVET) to explore advanced video coding technology, and the team developed Joint Exploration Model (JEM) for study. Moreover, since the year of 2018, the JVET team has been working on a new video coding standard, informally called Versatile Video Coding (VVC), as the successor of HEVC. It is anticipated that VVC may improve the compression efficiency by saving around 50% bits while maintaining the same quality, especially

for UHD video, compared to HEVC. Nonetheless, it is worth noting that the improvement of VVC is probably achieved at the cost of multiplicative encoding/decoding complexity.

### *B. Deep Learning for Image/Video Coding*

The past decade has witnessed the emerging and booming of deep learning, a class of techniques that are increasingly adopted in the hope of approaching the ultimate goal of artificial intelligence [59]. Deep learning belongs to machine learning technology, and has the distinction of its computational models, known as deep artificial neural networks or deep networks for short, which are composed of multiple (usually more than three) processing layers, each layer is further composed of multiple simple but non-linear basic computational units. One benefit of such deep networks is believed to be the capacity for processing data with multiple levels of abstraction, and converting data into different kinds of representations. Note that these representations are not manually designed; instead, the deep network including the processing layers is learned from massive data using a general machine learning procedure. Deep learning eliminates the necessity of handcrafted representations, and thus is regarded useful especially for processing natively unstructured data, such as acoustic and visual signal, whilst processing such data has been a longstanding difficulty in the artificial intelligence field.

Specifically for processing image/video, deep learning using convolutional neural network (CNN) has revolutionized the paradigm in computer vision and image processing. In 2012, Krizhevsky *et al.* [57] designed a 8-layer CNN, which won the image classification challenge by a surprisingly low error rate compared with previous works. In 2014, Girshick *et al.* [33] promoted the performance of object detection by a significant margin with the proposed regions with CNN features. Also in 2014, Dong *et al.* [27] proposed a 3-layer CNN for single image super-resolution (SR), which outperforms the previous methods in both reconstruction quality and computational speed. In 2017, Zhang *et al.* [142] presented a deep CNN for image denoising, and demonstrated that a single CNN model may tackle with several different image restoration tasks including denoising, single image SR, and compression artifact reduction, while these tasks had been studied separately for a long while.

Witnessing such successful cases, experts cannot help but ask whether deep learning can benefit image/video coding as well. In history, artificial neural network is not strange to the image/video coding community. From 1980s to 1990s, a number of researches were conducted on neural network-based image coding [28], [48], but then the networks were shallow and the compression efficiency was not satisfactory. Thanks to the abundance of data, the more and more powerful computing platform, and the development of advanced algorithms, it is now possible to train very deep networks with even more than 1000 layers [40]. Thus, the exploration of using deep learning for image/video coding is worthy of reconsideration, and indeed has been an actively developing research area since 2015. At present, researches have shown promising results, confirmed the feasibility of deep learning-based image/video coding. Nonetheless, this technology is far from mature and calls for much more research and development efforts.

In this paper, we aim to provide a comprehensive review of the newest reports about deep learning-based image/video coding (until the end of 2018), as well as to present a case study of our developed prototype video

codec namely Deep Learning Video Coding (DLVC), so as to make interested readers aware of the status quo. Readers may also refer to [84] for a recently published review paper about the same theme.

The remainder of this paper is organized as follows. Sections II and III provide a review of related works about using deep learning for image/video coding. The related works are divided into two categories, and reviewed in the two sections respectively. The first category is deep schemes, i.e. new coding schemes that are built primarily upon deep networks; the second category is deep tools, i.e. deep network-based coding tools that are embedded into traditional, non-deep coding schemes; a deep tool may either replace its counterpart in the traditional scheme, or be newly added into the scheme. Section IV presents the case study of our developed DLVC, with all the design details and experimental results. Section V summarizes our perspectives on some open problems for future research, and then concludes this paper. Table I lists the abbreviations used in this paper.

### C. Preliminaries

In this paper, we consider coding methods for natural image/video, i.e. the image/video as-is seen by human taken by daily cameras or mobile phones. Although the methods are usually generally applicable, they have been specifically designed for natural image/video, and may not perform that well for other kinds (e.g. biomedical, remote-sensing).

Currently, almost all the natural image/video is in digital format. A grayscale digital image can be denoted by  $\mathbf{x} \in \mathbb{D}^{m \times n}$ , where  $m$  and  $n$  are the number of rows (height) and number of columns (width) of the image, and  $\mathbb{D}$  is the definition domain of a single picture element (pixel). For example,  $\mathbb{D} = \{0, 1, \dots, 255\}$  is a common setting, where  $|\mathbb{D}| = 256 = 2^8$ , thus the pixel value can be represented by an 8-bit integer; accordingly, an uncompressed grayscale digital image has 8 bits-per-pixel (bpp), while compressed bits are definitely less.

A color image is usually decomposed into multiple channels to record the color information. For example, using the RGB color space, a color image can be denoted by  $\mathbf{x} \in \mathbb{D}^{m \times n \times 3}$ , where 3 corresponding to three channels—Red, Green, Blue. Since human vision is more sensitive to luminance than to chrominance, the YCbCr (YUV) color space is much more adopted than RGB, and the U and V channels are typically down-sampled to achieve compression. For example, in the so-called YUV420 color format, a color image can be denoted by  $\mathbf{X} = \{\mathbf{x}_Y \in \mathbb{D}^{m \times n}, \mathbf{x}_U \in \mathbb{D}^{\frac{m}{2} \times \frac{n}{2}}, \mathbf{x}_V \in \mathbb{D}^{\frac{m}{2} \times \frac{n}{2}}\}$ .

A color video is composed by multiple color images, called frames, to record the scene at different timestamps. For example, in the YUV420 color format, a color video can be denoted by  $\mathbf{V} = \{\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_{T-1}\}$  where  $T$  is the number of frames, each  $\mathbf{X}_i = \{\mathbf{x}_Y^{(i)} \in \mathbb{D}^{m \times n}, \mathbf{x}_U^{(i)} \in \mathbb{D}^{\frac{m}{2} \times \frac{n}{2}}, \mathbf{x}_V^{(i)} \in \mathbb{D}^{\frac{m}{2} \times \frac{n}{2}}\}$ . If  $m = 1080, n = 1920, |\mathbb{D}| = 2^{10}$ , and a video has 50 frames-per-second (fps), then the data rate of the uncompressed video is  $1080 \times 1920 \times (10 + \frac{10}{4} + \frac{10}{4}) \times 50 = 1,555,200,000$  bits-per-second (bps), about 1.555 Gbps. Obviously, the video should be compressed by a ratio of hundreds to thousands before it can be efficiently transmitted over the current wired and wireless networks.

The existing lossless coding methods can achieve a compression ratio of about 1.5 to 3 for natural image, which is clearly below requirement. Thus, lossy coding is introduced to compress more but at the cost of incurring loss. The

TABLE I  
LIST OF ABBREVIATIONS

Abbreviation	Remark
AVC	Advanced Video Coding, i.e. H.264 [128]
BARC	block adaptive resolution coding
BD-rate	Bjontegaard's delta-rate [13]
BPG	Better Portable Graphics, an image coding format based on HEVC
CNN	convolutional neural network
CTU	coding tree unit
CU	coding unit
DLVC	Deep Learning Video Coding, our developed prototype video codec
GAN	generative adversarial network
HEVC	High Efficiency Video Coding, i.e. H.265 [108]
HM	HEVC reference software
ILF	in-loop filter
JEM	Joint Exploration Model, a video coding software developed by JVET
JPEG <sup>1</sup>	Joint Photographic Experts Group, a group of ISO/IEC
JPEG <sup>2</sup>	a standard published by ISO/IEC [121]
JVET	Joint Video Experts Team, a team of MPEG and VCEG
LSTM	long short-term memory
MAE	mean-absolute-error
MC	motion compensation
ME	motion estimation
MPEG	Moving Picture Experts Group, a group of ISO/IEC
MSE	mean-squared-error
MS-SSIM	multi-scale SSIM
PSNR	peak signal-to-noise ratio
QP	quantization parameter
ReLU	rectified linear unit [90]
RNN	recurrent neural network
SR	super-resolution
SSIM	structural similarity [126]
VCEG	Video Coding Experts Group, a group of ITU-T
VVC	Versatile Video Coding, an incoming video coding standard

loss can be measured by the difference between original and reconstructed images, e.g. using mean-squared-error (MSE) for grayscale image:

$$\text{MSE} = \frac{\|x - x_{\text{rec}}\|^2}{m \times n} \quad (1)$$

Accordingly, the quality of reconstructed image compared with original image can be measured by peak signal-to-noise ratio (PSNR):

$$\text{PSNR} = 10 \times \log_{10} \frac{(\max(\mathbb{D}))^2}{\text{MSE}} \quad (2)$$

where  $\max(\mathbb{D})$  is the maximal value in  $\mathbb{D}$ , e.g. 255 for 8-bit grayscale image. For color image/video, the PSNR values of Y, U, V are usually separately calculated. For video, the PSNR values of different frames are usually

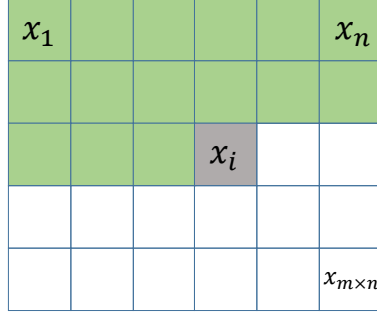


Fig. 1. Illustration of a typical predictive coding scheme, where the pixels are encoded/decoded one by one in the raster scan order. For the pixel  $x_i$  (marked gray), all the previous pixels (marked green), i.e. above of  $x_i$  and left in the same row of  $x_i$ , can be used as condition to predict the pixel value of  $x_i$ . The green area is also called the *context* for  $x_i$ . For simplification, context can be chosen as a subset of the green area.

separately calculated and then averaged. There are other quality metrics in replacement of PSNR, such as structural similarity (SSIM) and multi-scale SSIM (MS-SSIM) [126].

To compare different lossless coding schemes, it is sufficient to compare the compression ratio, or the resulting rate (bpp, bps, etc.). To compare different lossy coding schemes, it is necessary to take into account both rate and quality. For example, to calculate the relative rates at several different quality levels, and then to average the rates, is a commonly adopted method; the average relative rate is known as Bjontegaard's delta-rate (BD-rate) [13]. There are other important aspects to evaluate image/video coding schemes, including encoding/decoding complexity, scalability, robustness, and so on.

## II. REVIEW OF DEEP SCHEMES

In this section we review some representative coding schemes that are built primarily upon deep networks. Generally speaking, there are two approaches for deep image coding schemes, i.e. pixel probability modeling and auto-encoder. These two approaches are combined together in several deep schemes. In addition, we discuss deep video coding schemes and special-purpose coding schemes, where special-purpose schemes are further categorized into perceptual coding and semantic coding.

### A. Pixel Probability Modeling

According to Shannon's information theory [102], the optimal method for lossless coding can reach the minimal coding rate  $-\log_2 p(x)$  where  $p(x)$  is the probability of the symbol  $x$ . To reach this target, a number of lossless coding methods have been invented, and arithmetic coding is believed to be among the optimal ones [129]. In essence, given a probability  $p(x)$ , arithmetic coding ensures the coding rate to be as near as possible to  $-\log_2 p(x)$  up to rounding error. Thus, the remaining problem is to find out the probability, which is however very difficult for natural image/video as it is of very high dimension.

One way to estimate  $p(\mathbf{x})$ , where  $\mathbf{x}$  is an image, is to decompose the image into  $m \times n$  pixels and to estimate the probabilities of these pixels one by one (e.g. in the raster scan order). This is a typical *predictive coding* strategy. Note that

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \cdots p(x_i|x_1, \dots, x_{i-1}) \cdots p(x_{m \times n}|x_1, \dots, x_{m \times n-1}) \quad (3)$$

which is illustrated in Fig. 1. Here the condition for  $x_i$  is also called the *context* for  $x_i$ . When the image is large, the conditional probability can be difficult to estimate. A simplification is to reduce the range of context, e.g.

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \cdots p(x_i|x_{i-k}, \dots, x_{i-1}) \cdots p(x_{m \times n}|x_{m \times n-k}, \dots, x_{m \times n-1}) \quad (4)$$

where  $k$  is a prechosen constant.

As known, deep learning is good at solving regression and classification problems. Therefore, it has been proposed to estimate the probability  $p(x_i|x_1, \dots, x_{i-1})$  given the context  $x_1, \dots, x_{i-1}$ , using trained deep networks. This strategy is proposed for other kinds of high-dimensional data in as early as 2000 [12], but is applied to image/video until recently. For example, in [58], the probability estimation is considered for binary images, i.e.  $x_i \in \{-1, +1\}$ , where it suffices to predict a single probability value  $p(x_i = +1|x_1, \dots, x_{i-1})$  for each pixel. The paper presents the neural autoregressive distribution estimator (NADE), where a feed-forward network with one hidden layer is used for each pixel, and the parameters are shared across these networks. The parameter sharing also help to speed up the computations for each pixel. A similar work is presented in [37], where the feed-forward network also has connections skipping the hidden layer, and the parameters are also shared. Both [58] and [37] perform experiments on the binarized MNIST dataset<sup>1</sup>. Uria *et al.* [116] extend the NADE to real-valued NADE (RNADE), where the probability  $p(x_i|x_1, \dots, x_{i-1})$  is made up by a mixture of Gaussians, and the feed-forward network needs to output a set of parameters for the Gaussian mixture model, instead of a single value in NADE. Their feed-forward network has one hidden layer and parameter sharing, but the hidden layer is equipped with rescaling to avoid saturation, and uses rectified linear unit (ReLU) [90] instead of sigmoid. They also consider mixture of Laplacians rather than Gaussians. Experiments are conducted on  $8 \times 8$  natural image patches, where the pixel value is added with noise and converted to real value. In [117], NADE and RNADE are improved by using different orderings of the pixels as well as using more hidden layers in the network. In [120], RNADE is improved by enhancing the Gaussian mixture model (GMM) with deep GMM.

Designing advanced networks has been an important theme for improving pixel probability modeling. In [109], multi-dimensional long short-term memory (LSTM) based network is proposed, together with mixtures of conditional Gaussian scale mixtures, a generalization of GMM, for probability modeling. LSTM is a kind of recurrent neural networks (RNNs), and is regarded good at modeling sequential data. The spatial variant of LSTM is used for images. Later in [118], several different networks are studied, including RNNs and CNNs that are known as PixelRNN and PixelCNN, respectively. For PixelRNN, two variants of LSTM, called row LSTM and diagonal BiLSTM, are proposed, where the latter is specifically designed for images. PixelRNN incorporates residual connections [40] to help train deep networks with up to 12 layers. For PixelCNN, in order to suit for the shape of context (see Fig.

<sup>1</sup>The raw MNIST dataset can be accessed at <http://yann.lecun.com/exdb/mnist/>.

1), *masked* convolutions are proposed. PixelCNN is also as deep as having 15 layers. Compared with previous works, PixelRNN and PixelCNN are more dedicated to natural images: they consider pixels as discrete values (e.g.  $0, 1, \dots, 255$ ), and predict a multinomial distribution over the discrete values; they deal with color images (in RGB color space); multi-scale PixelRNN is proposed; and they work well on the CIFAR-10 and ImageNet datasets. Quite a number of researches follow the approach of PixelRNN and PixelCNN. In [119], Gated PixelCNN is proposed to improve the PixelCNN, and achieves comparable performance with PixelRNN but with much less complexity. In [99], PixelCNN++ is proposed with the following improvements upon PixelCNN: a discretized logistic mixture likelihood is used rather than a 256-way multinomial distribution; down-sampling is used to capture structures at multiple resolutions; additional short-cut connections are introduced to speed up training; dropout is adopted for regularization; RGB is combined for one pixel. In [18], PixelSNAIL is proposed, in which casual convolutions are combined with self attention.

Most of the aforementioned works directly model pixel probability. In addition, pixel probability may be modeled as a conditional one upon explicit or latent representations. That says, we may estimate

$$p(\mathbf{x}|\mathbf{h}) = \prod_{i=1}^{m \times n} p(x_i|x_1, \dots, x_{i-1}, \mathbf{h}) \quad (5)$$

where  $\mathbf{h}$  is the additional condition. Note also that  $p(\mathbf{x}) = p(\mathbf{h})p(\mathbf{x}|\mathbf{h})$ , which means the modeling is split into an unconditional and a conditional. For example, in [119], the additional condition can be image class or high-level image representations that are derived by another deep network. In [56], PixelCNN with latent variables are considered, where the latent variables are derived from the original image: they can be a quantized grayscale version of the original color image, or a multi-resolution image pyramid.

Regarding practical image coding schemes, in [64], a network with *trimmed* convolutions is adopted to predict probabilities for binary data, while a 8-bit grayscale image with the size of  $m \times n$  is converted to a binary cube with the size of  $m \times n \times 8$  to be processed by the network. The network is similar to PixelCNN but is of three dimension. The trimmed convolutional network-based arithmetic encoding (TCAE) is reportedly better than the previous non-deep lossless coding schemes, such as TIFF, GIF, PNG, JPEG-LS, and JPEG 2000-LS; on the Kodak image set<sup>2</sup>, TCAE achieves a compression ratio of 2.00. Differently, in [4], CNN is used in the wavelet transform domain rather than the pixel domain, i.e. CNN is to predict wavelet detail coefficients from coefficients within neighboring subbands.

For video coding, in [52], PixelCNN is generalized to video pixel network (VPN) for the pixel probability modeling of video. VPN is composed of CNN encoders (for previous frames to predict the current frame) and PixelCNN decoders (for prediction inside the current frame). CNN encoders preserve at all layers the spatial resolution of input frames to maximize representational capacity. Dilated convolutions are adopted to enlarge receptive fields and better capture global motion. The outputs of the CNN encoders are combined over time with a convolutional LSTM that also preserves the resolution. The PixelCNN decoders use masked convolutions and

<sup>2</sup>The Kodak image set can be accessed at <http://www.r0k.us/graphics/kodak/>.



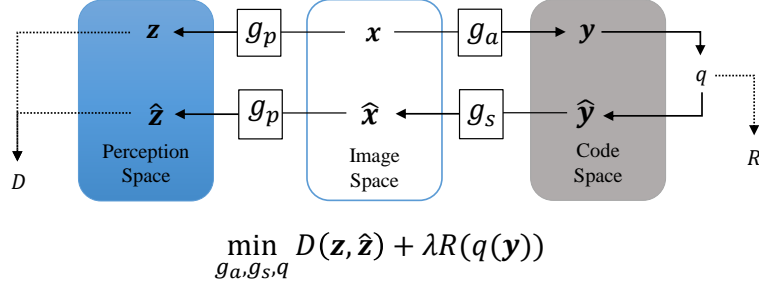


Fig. 2. Illustration of a typical transform coding scheme. The original image  $\mathbf{x}$  is transformed by an *analysis* function  $g_a$  to achieve the code  $\mathbf{y}$ . The code  $\mathbf{y}$  is quantized (denoted by  $q$ ) and compressed into bits. The number of bits is used to measure the coding rate ( $R$ ). The quantized code  $\hat{\mathbf{y}}$  is then inversely transformed by a *synthesis* function  $g_s$  to achieve the reconstructed image  $\hat{\mathbf{x}}$ . Both of  $\mathbf{x}$  and  $\hat{\mathbf{x}}$  are further transformed by a same *perceptual* function  $g_p$ , resulting in  $\mathbf{z}$  and  $\hat{\mathbf{z}}$ , respectively. The difference between  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  is used to measure the distortion ( $D$ ).

adopt multinomial distributions over discrete pixel values. VPN is experimented on the Moving MNIST and Robotic Pushing datasets.

In addition, Schioppa *et al.* [101] investigate a lossless image coding scheme, where they use CNN to predict pixel value rather than its distribution. The predicted value is subtracted from the true pixel value, resulting in residue that is then coded. In addition, they consider the adaptive selection among the CNN predictor and some non-CNN predictors.

### B. Auto-Encoder

Auto-encoder originates from the well-known work of Hinton and Salakhutdinov [42], which trains a network for dimensionality reduction and the network consists of encoding part and decoding part. The encoding part converts an input high-dimension signal to its low-dimension representation, and the decoding part recovers (not perfectly) the high-dimension signal from the low-dimension representation. Auto-encoder enables automated learning of representations and eliminates the need of hand-crafted features, which is also believed to be one of the most important advantages of deep learning.

It seems quite straightforward to adopt the auto-encoder network for lossy image coding: the encoder and decoder are trained out, and we just need to encode the learned representation. However, the traditional auto-encoder is not optimized for compression, and directly using a trained auto-encoder is not efficient [127]. When we consider the compression requirement, there are several challenges: First, the low-dimension representation shall be quantized then coded, but the quantization step is not differentiable, making a difficulty to train the network. Second, lossy coding is to achieve a better tradeoff between rate and quality, so the rate shall be taken into account when training the network, but the rate is not easy to calculate or estimate. Third, a practical image coding scheme needs to consider variable rate, scalability, encoding/decoding speed, interoperability, and so on. In response to these challenges, a number of researches have been conducted especially in recent years.

A conceptual illustration of auto-encoder-based image coding scheme is shown in Fig. 2, which is a typical *transform coding* strategy. The original image  $\mathbf{x}$  is transformed to  $\mathbf{y} = g_a(\mathbf{x})$ , and  $\mathbf{y}$  is quantized then coded. The

decoded  $\hat{y}$  is inversely transformed to  $\hat{x} = g_s(\hat{y})$ . Considering the tradeoff between rate and quality, we can train the network to minimize the joint rate-distortion cost  $D + \lambda R$  where  $D$  is calculated or estimated as the difference between  $x$  and  $\hat{x}$  (note that the difference may be calculated or estimated in a perception space),  $R$  is calculated or estimated from the quantized code, and  $\lambda$  is the Lagrange multiplier. All of the existing researches follow this scheme more or less and differ in their network structure and loss function.

For the network structure, RNNs and CNNs are the widely used two categories. The most representative works include:

- Toderici *et al.* [111] propose a general framework for variable rate image compression. They use binary quantization to generate codes, and do not consider the rate during training, i.e. the loss is only end-to-end distortion, measured by MSE. Their framework indeed provides a scalable coding functionality, where RNN (specifically LSTM) with convolutional and deconvolutional layers is reported to perform well. They provide results on a large-scale dataset of  $32 \times 32$  thumbnails. Later, Toderici *et al.* [112] propose an improved version, where they use a neural network like PixelRNN [118] to compress the binary codes; they also introduce a new gated recurrent unit (GRU) inspired by the residual network (ResNet) [40]. They report better results than JPEG on the Kodak image set using MS-SSIM as quality metric. Johnston *et al.* [51] further improve the RNN-based method by introducing hidden-state priming into RNN, using an SSIM-weighted loss function, and enabling spatially adaptive bitrates. They achieve better results than BPG on the Kodak image set using MS-SSIM. Covell *et al.* [22] enable spatially adaptive bitrates by training stop-code tolerant RNNs.
- Ballé *et al.* [9] propose a general framework for rate-distortion optimized image compression. They use multiary quantization to generate integer codes and consider the rate during training, i.e. the loss is the joint rate-distortion cost, where distortion can be MSE or others. To estimate the rate, they use adding a random noise to replace the quantization during training, and use the differential entropy of the noisy “code” as a proxy for the rate. As for the network structure, they use the generalized divisive normalization (GDN) transform, which consists of a linear mapping (matrix multiplication) followed by a nonlinear parametric normalization; the effectiveness of the proposed GDN for image coding is verified in [8]. Later, Ballé *et al.* [10] propose an improved version, where they use 3 convolutional layers each followed by down-sampling and a GDN operation to implement the transform; accordingly, they use 3 layers of inverse GDN + up-sampling + convolution to implement the inverse transform. In addition, they design an arithmetic coding method to compress the integer codes. They report better results than JPEG and JPEG 2000 on the Kodak image set using MSE as quality metric. Furthermore, Ballé *et al.* [11] improve their scheme by incorporating a scale hyper-prior into the auto-encoder, which is inspired by the variational auto-encoder [55]. They use another transform  $h_a$  to convert  $y$  into  $w = h_a(y)$ , quantize and encode  $w$  (transmitted as side information), and use another inverse transform  $h_s$  to convert the decoded  $\hat{w}$  into the estimated standard deviation of the quantized  $\hat{y}$ , which is then used during the arithmetic coding of  $\hat{y}$ . On the Kodak image set and using PSNR as quality metric, their method is only slightly worse than BPG.

Besides [9], several works also concentrate on dealing with the non-differentiable quantization and/or the es-

timization of rate. Theis *et al.* [110] adopt a very simple work-around for quantization: quantization is performed as usual in the forward pass, but the gradients are directly passed through the quantization layer in the backward pass. Surprisingly this work-around works well. In addition, they replace the rate with an upper bound that is differentiable. Dumas *et al.* [29] consider a stochastic winner-take-all mechanism, where the entries in  $\mathbf{y}$  with the largest absolute values are kept and the other entries are set to 0; then the entries are uniformly quantized and compressed. Agustsson *et al.* [2] propose a soft-to-hard vector quantization scheme, where they use a soft quantization (i.e. assigning a representation to multiple codes with different membership values) rather than hard quantization (i.e. assigning a representation to only one code) during training, and they adopt an annealing process to let the soft quantization approach the hard quantization gradually. Note that their scheme takes advantage of vector quantization while other works usually adopt scalar quantization. Li *et al.* [65] introduce an importance map for rate estimation, where the importance map is quantized to a mask and the mask decides how many bits are kept at each location, thus the sum of the importance map can be used as a rough estimate of the coding rate.

Besides [111], several works also consider the functionality of variable rate with less or no training for different rates. In [110], scale parameters are introduced and a pretrained auto-encoder is fine-tuned for different rates. In [30], a unique learned transform is proposed, together with variable quantization step for different rates. In [15], a multi-scale decomposition transform is trained and optimized for all scales; and rate allocation algorithms are provided to determine the optimal scale of each image block for either a target rate or a target quality factor. Besides, scalable coding is considered in [146] differently from that in [111]. In [146], an image is decomposed into multiple bit-planes, which are transformed and quantized in parallel; bidirectional assembling gated units are proposed to reduce the correlation between different bit-planes.

Several works consider advanced network structures and different loss functions. Theis *et al.* [110] adopt a sub-pixel structure for computational efficiency. Rippel and Bourdev [97] present a pyramid decomposition followed by inter-scale alignment network, which is lightweight and runs in real-time. They also use a discriminator loss in addition to the reconstruction loss. Snell *et al.* [104] use the MS-SSIM as loss function instead of MSE or mean-absolute-error (MAE) to train auto-encoders, and they find that MS-SSIM is better calibrated to perceptual quality. Zhou *et al.* [149] use deeper networks for encoder/decoder and a separate network for post-processing at the decoder side. They also replace the Gaussian model in [11] with the Laplacian model.

As mentioned before, pixel probability modeling represents predictive coding and auto-encoder represents transform coding. These two strategies can be combined for higher compression efficiency. Mentzer *et al.* [87] propose a practical lossless image coding scheme, where they use auto-encoders at multiple levels to learn the condition for pixel probability modeling. Mentzer *et al.* [86] integrate pixel probability modeling (a 3D PixelCNN) into auto-encoder so as to estimate the coding rate and to train the PixelCNN and the auto-encoder jointly. Baig *et al.* [6] introduce partial-context image inpainting into the variable rate compression framework [111], which is actually to predict a block from the block's context, assuming the blocks are encoded/decoded one by one in the raster scan order (similar to what is shown in Fig. 1 but at the block level). The prediction signal is added onto the network output signal to achieve  $\hat{x}$ , i.e. the transform coding network deals with the prediction residues. Minnen *et al.* [89] additionally consider rate allocation among the blocks. Similarly but in a different manner, Minnen *et*

*al.* [88] improve upon [11] by augmenting the hyper-prior with the context, i.e. they use not only  $\hat{w}$  but also the context to predict the probability of each entry of  $\hat{y}$ . Their method outperforms BPG on the Kodak image set and using PSNR as quality metric, which represents the state of the art by the end of 2018. Lee *et al.* [60] introduce the context adaptive entropy model into the hyper-prior  $\hat{w}$ .

Moreover, Cheng *et al.* [21] apply principle component analysis on the learned representation, which is virtually a second transform.

### C. Video Coding

Starting from 2017, a few researches have been reported for deep video coding schemes. Compared to image coding, video coding calls for efficient methods to remove the inter-picture redundancy. Inter-picture prediction is then an important issue in these researches. Motion estimation and compensation is widely adopted, but is implemented by trained deep networks until recently.

Chen *et al.* [17] seems the first to report a video coding scheme by using trained deep networks as auto-encoders. Specifically, they divide video frames into  $32 \times 32$  blocks and for each block they choose one from two modes: intra coding or inter coding. If using intra coding, there is an auto-encoder to compress the block. If using inter coding, then they perform motion estimation and compensation using the traditional method, and input the residues to another auto-encoder. For both auto-encoders, the encoded representations are directly quantized and coded by the Huffman method. This scheme is quite rough and does not compete H.264.

Wu *et al.* [131] propose a video coding scheme with image interpolation, where the key frames (I frames) are first compressed by the deep image coding scheme in [112], and the remaining frames (B frames) are then compressed in a hierarchical order. For each B frame, two compressed frames (either I frames or previously compressed B frames) before and after are used to “interpolate” the current frame: the motion information is used to warp the two compressed frames (i.e. motion compensation), and then the two warped frames are sent as side information to a variable rate image coding scheme that processes the current frame. The scheme is reported to perform on par with H.264.

Chen *et al.* [20] propose another video coding scheme with the so-called PixelMotionCNN. In their scheme, frames are compressed in the temporal order, and each frame is divided into blocks that are compressed in the raster scan order. Before one frame is compressed, the previous two compressed frames are used to “extrapolate” the current frame. When a block is to be compressed, the extrapolated frame together with the block’s context are sent to the PixelMotionCNN to generate a prediction signal for the current block, then the prediction residues are compressed by the variable rate image coding scheme in [112]. This scheme also performs on par with H.264.

Lu *et al.* [80] propose a real end-to-end deep video coding scheme, which can be viewed as a “deepened” version of the traditional video coding schemes. Specifically in their scheme, for each frame to be compressed, an optical flow estimation module is used to obtain the motion information between the frame and the previous compressed frames. Motion compensation is also performed by a trained network, to generate a prediction signal for the current frame. For the prediction residues and the motion information, two auto-encoders are used to compress them respectively. The entire network is jointly optimized with a single loss function, i.e. the joint rate-distortion

cost. This scheme reportedly achieves better compression efficiency than H.264, and even outperforms HEVC (x265 encoder) when evaluated with MS-SSIM.

Rippel *et al.* [98] present the to-date most sophisticated deep video coding scheme, which inherits and extends a deepened version of the traditional video coding schemes. Their scheme has the following new features: (1) only one auto-encoder to compress motion information and prediction residues simultaneously; (2) a state that is learned from the previous frames and updated recursively; (3) motion compensation with multiple frames and multiple optical flows; (4) a rate control algorithm. This scheme is reported to outperform HEVC reference software (HM) when evaluated with MS-SSIM.

By the end of 2018, we do not observe any report that a deep video coding scheme can outperform HM when evaluated with PSNR, which seems a hard mission.

#### D. Special-Purpose Coding

Most of the researches about deep schemes concern image/video coding for *signal fidelity*, i.e. to minimize the distortion between original and reconstructed image/video subject to a given rate, where the distortion can be defined as MSE or other differences. However, if we do not concern the fidelity, we may instead care about the perceptual naturalness of the reconstructed image/video, or the utility of the reconstructed image/video in semantic analysis tasks. The latter two kinds of quality metrics are termed *perceptual naturalness* and *semantic quality*. There have been a few works that tailor image/video coding for these quality metrics.

1) *Perceptual Coding*: Since the boom of generative adversarial network (GAN) [34], deep networks are known to be capable in generating perceptually natural images. Leveraging this capability at the decoder side can surely improve the perceptual quality of decoded images. Different from the generator in normal GANs, the decoder should also ensure the decoded images to be similar to original images, which raises a problem of controlled generation and the encoder actually provides the control signal in the coded bits.

Inspired by the variational auto-encoder (VAE) [55], Gregor *et al.* [36] propose Deep Recurrent Attentive Writer (DRAW) for image generation, which extends the traditional VAE by using RNNs as encoder and decoder. Unfolding the encoder RNN produces a series of latent representations. Then, Gregor *et al.* [35] introduce convolutional DRAW, and observe that it is able to transform an image into a series of increasingly detailed representations, ranging from global conceptual aspects to low level details. Thus, they suggest a conceptual compression scheme, whose one benefit is to achieve plausible reconstruction images at very low bit rates.

It has been realized that perceptual naturalness can be evaluated by the discriminator in GAN [14]. Several researches are devoted to deep coding schemes for perceptual quality using the discriminator loss solely or jointly with MSE or other losses. For example, Santurkar *et al.* [100] propose the so-called generative compression schemes for both image and video. For image, they first train a canonical GAN, then they use the generator as the decoder, fix it, and train the encoder to minimize a sum of MSE and feature loss. For video, they reuse the encoder and decoder trained for image, transmit only a few frames, and restore the other frames at the decoder side via interpolation. Their schemes are able to achieve very high compression ratio. Kim *et al.* [54] build a new video compression scheme, where a few key frames are normally compressed (by H.264) and the other frames are extremely compressed.

Indeed, edges are extracted from the down-sampled non-key frames and transmitted. At the decoder side, the key frames are firstly reconstructed, then edges are similarly extracted from them. A conditional GAN is trained with the reconstructed key frames where edge is the condition. Then the conditional GAN is used to generate the non-key frames. Again, their scheme performs well at very low bit rates.

2) *Semantic Coding*: A few researches have been conducted on deep coding schemes that preserve the semantic information or concern the semantic quality.

Agustsson *et al.* [3] present a GAN-based image compression scheme for extremely low bit rates. The scheme combines auto-encoder and GAN, collapsing the decoder and the generator into one. In addition, a semantic label map can be used as an additional input to the encoder, and as a condition for the discriminator. It is reported that the proposed scheme reconstructs images with higher semantic quality, in the sense that the semantic segmentation on these images is more accurate than that on BPG-compressed images at the same rate.

Luo *et al.* [81] propose a concept of deep semantic image compression (DeepSIC), which incorporates the semantic information (e.g. classes) into the coded bits. There are two versions of DeepSIC, both based on auto-encoder. In the one version, the semantic information is extracted from the representation  $y$ , and encoded into the bits. In the other version, the semantic information is not encoded, but extracted at the decoder side from the quantized representation  $\hat{y}$ . Torfason *et al.* [113] investigate performing semantic analysis tasks (classification and semantic segmentation) from the quantized representation rather than from the reconstructed image. That says, the decoding process is omitted. They show that the classification and segmentation accuracy values are very close between the representation and the image, but the computational complexity is reduced significantly. Zhang *et al.* [143] study a deep image coding scheme for simultaneous compression and retrieval. Their motivation is that the coded bits can be used not only for reconstructing image but also for retrieving similar images *without decoding*. They use an auto-encoder to compress image into bits, and use a revised classification network to extract binary features. Then they combine the two parts of bits, and fine-tune the feature extraction network for image retrieval. Their results indicate that at the same rate, the reconstructed images are better than JPEG-compressed ones, and the retrieval accuracy improves due to the fine-tuning.

Akbari *et al.* [5] design a scalable image coding scheme where the coded bits consist of three layers. The first layer is the semantic segmentation map coded losslessly. The second layer is a down-sampled version of the original image also coded losslessly. With the first two layers, a network is trained to predict the original image and the prediction residues are coded by BPG as the third layer. This scheme is reported to outperform BPG on the Kodak image set when evaluated with PSNR and MS-SSIM.

Chen and He [19] consider deep coding for facial images with semantic quality metric instead of PSNR or perceptual quality. For this purpose, their loss function has three parts: MAE, discriminator loss, and a semantic loss, where the semantic loss is to project the original and reconstructed images into a compact Euclidean space through a learned transformation, and to calculate the Euclidean distance between them. Accordingly, their scheme performs very well when evaluated with face verification accuracy at the same rate.

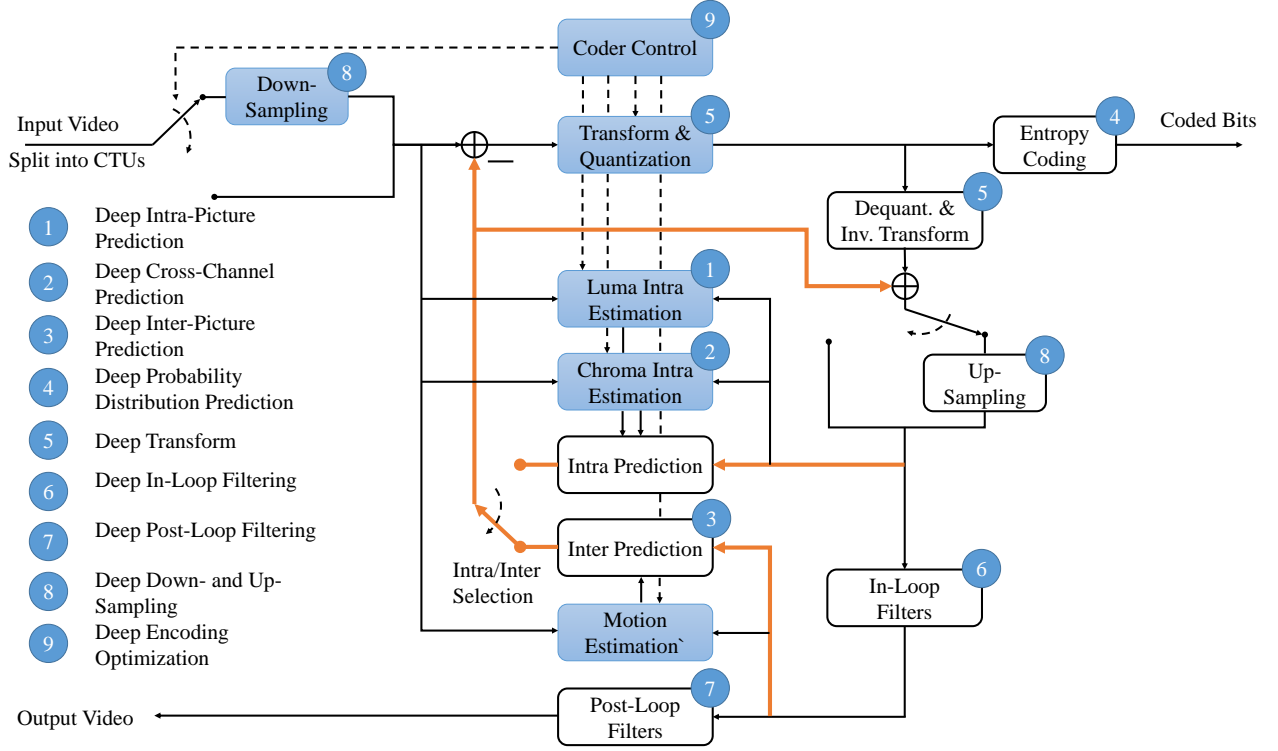


Fig. 3. Illustration of a traditional hybrid video coding scheme as well as the locations of deep tools inside the scheme. Note that the yellow lines indicate the flow of prediction, and the blue boxes indicate the tools that are used at the encoder side only.

### III. REVIEW OF DEEP TOOLS

In this section we review some representative works about using trained deep networks as tools within the traditional coding schemes or together with traditional coding tools. Generally speaking, the traditional video coding schemes adopt a hybrid coding strategy, i.e. a combination of predictive coding and transform coding. As depicted in Fig. 3, an input video sequence is divided into pictures, pictures are divided into blocks (the largest block is called CTU, which can be divided into smaller CUs, in HEVC [108]), and blocks are divided into channels (i.e. Y, U, V). The pictures/blocks/channels are compressed in a predefined order, and the previously compressed ones can be used to predict the following ones, which is known as intra-picture prediction (between blocks), cross-channel prediction (between channels), and inter-picture prediction (between pictures), respectively. The prediction residues are then transformed and quantized and entropy coded to achieve the final bits. Some auxiliary information such as block partition and prediction mode is also entropy coded into the bits (not shown in the figure). Probability distribution prediction is used in the entropy coding step. Since the quantization step loses information and may cause artifacts, filtering is proposed to enhance the reconstructed video, which may be performed in-loop (before predicting the next picture) or out-of-loop (before output). In addition, to reduce the data volume, the pictures/blocks/channels may be down-sampled before being compressed, and up-sampled afterwards. Finally, the encoder needs to control the different modules and combine them to achieve a tradeoff between coding rate, quality, and computational speed. Encoding optimization is an important theme in practical coding systems.

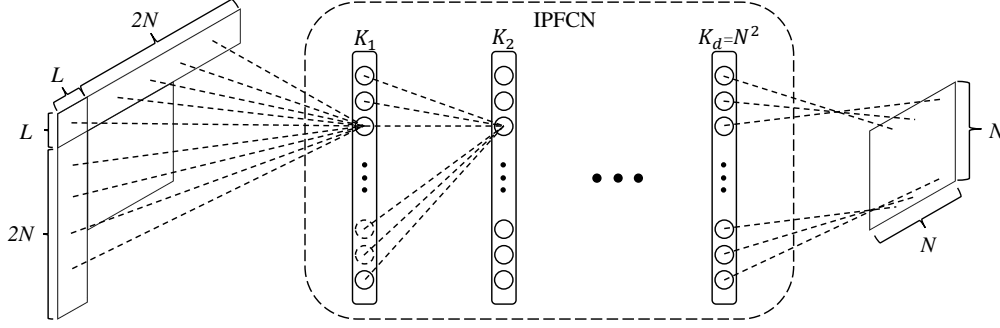


Fig. 4. Illustration of a fully connected network for intra prediction (IPFCN).

Trained deep networks can act as almost all of the modules shown in Fig. 3, where we have indicated different locations for deep tools. In the following, we will review the works about deep tools according to where they are used in the entire scheme.

#### A. Intra-Picture Prediction

Intra-picture prediction, or intra prediction for short, is a tool to predict between blocks inside the same picture. H.264 introduces intra prediction with several predefined prediction modes, such as DC prediction and extrapolation along different directions [128]. The encoder can choose a prediction mode for each block and signal the choice to the decoder. To decide mode, it is a common strategy to compare the coding rate and distortion of different modes and to select the mode with the minimal rate-distortion cost. In HEVC, more prediction modes are introduced [108].

Li *et al.* [63] propose a fully connected network for intra prediction that is depicted in Fig. 4. For the current  $N \times N$  block, they use  $L$  rows above and  $L$  columns to the left, in total  $4NL + L^2$  pixels as context. They use an image set known as the New York City Library to generate training data, in which the raw image is compressed at different quantization parameters (QPs). When training the network, they investigate two strategies: the first is to train a single model with all training data, and the second is to split the training data into two groups by considering the HEVC prediction modes, and to train two models respectively. The strategy of two models turns out better for compression. They integrate the trained networks as new prediction modes along with the HEVC modes. They report around 3% BD-rate reduction than HM.

Pfaff *et al.* [94] also adopt fully connected network for intra prediction, but propose to train multiple networks as different prediction modes. In addition, they propose to train a separate network whose input is also the block's context but output is the predicted likelihood of different modes. Moreover, they propose to use a different transform for each of the network-based prediction modes. Their reported performance is high: around 6% BD-rate reduction than an improved version of HM (with advanced block partitioning).

Hu *et al.* [44] devise a progressive spatial RNN for intra prediction. Different from the above works, they propose to leverage the sequential modeling capacity of RNN to generate prediction progressively from the context to the block. In addition, they suggest the use of sum-of-absolute-transformed-difference (SATD) as the loss function and argue that SATD correlates better to the rate-distortion cost.



Cui *et al.* [23] consider a CNN for intra prediction, or more specifically, intra prediction refinement. They use the HEVC prediction modes to generate prediction, and then use a trained CNN to refine the prediction. Note that the CNN has not only the HEVC prediction but also the context as its input. This method seems achieving only marginal gain.

### B. Inter-Picture Prediction

Inter-picture prediction, or inter prediction for short, is a tool to predict between video frames so as to remove the redundancy along the temporal dimension. Inter prediction is the kernel of video coding and it largely decides the compression efficiency of a video coding scheme. In the traditional video coding schemes, inter prediction is mostly fulfilled by block-level motion estimation (ME) and motion compensation (MC). Given a reference frame and a block to be coded, ME is to find the location in the reference frame where the content is the most similar to that inside the to-be-coded block, and MC is to retrieve the content at the found location so as to predict the block. Many techniques have been proposed to improve block-level ME and MC, such as using multiple reference frames, bi-directional inter prediction (i.e. using two reference frames jointly), fractional-pixel ME and MC, and so on.

Inspired by the multiple reference frames, Lin *et al.* [71] propose a new inter prediction mechanism by extrapolating the multiple reference frames. Specifically they adopt a Laplacian pyramid of GANs to extrapolate a frame from the previously compressed four frames. This extrapolated frame serves as another reference frame. They report around 2% BD-rate reduction than HM.

Inspired by the bi-directional inter prediction, Zhao *et al.* [148] propose a method to enhance the prediction quality. The previous bi-directional prediction simply computes a linear combination of two prediction blocks. They propose to employ trained CNN to combine the two prediction blocks in a nonlinear and data-driven manner.

Inspired by the fractional-pixel ME and MC, a number of researches are conducted on the fractional-pixel interpolation problem, which aims at generating imaginary pixels at fractional locations on the reference frame because the motion between two frames is not aligned to integer pixels. Here, a major difficulty is how to prepare training data because fractional pixels are imaginary. Yan *et al.* [137] propose to use a CNN for half-pixel interpolation, where they suggest a method that blurs a high-resolution image and then samples pixels from the blurred image: odd locations as integer pixels and even locations as half pixels. This method is inherited in [76], where the authors analyze the effect of different blurring degrees. Zhang *et al.* [141] propose another method, which formulates the fractional interpolation as a resolution enhancement problem. Thus, they down-sample high-resolution images to achieve training data. Yan *et al.* [136] consider a different formulation, treating the fractional-pixel MC as an inter-picture regression problem. They use video sequences to retrieve training data, where they rely on the fractional-pixel ME to align different frames, and use reference frame as integer pixels and current frame as fractional pixels. Yan *et al.* [135] further discover a key characteristic of the fractional interpolation problem, namely its invertibility: if fractional pixels can be interpolated from integer pixels, then integer pixels should also be interpolated from fractional pixels. Based on the invertibility, they propose an unsupervised manner to train CNN for half-pixel interpolation.

In addition to the improvements of inter prediction methods, another approach is considered where intra and inter predictions are combined. Specifically, the generation of prediction signal is based on not only reference frame but also context in current frame. For example, Huo *et al.* [45] propose to use a trained CNN to refine the inter prediction signal. They find that using the context of the to-be-predicted block can improve the prediction quality. Similarly, Wang *et al.* [124] also refine the inter prediction signal by a CNN, where the CNN inputs include the inter prediction signal, the context of the current block, and the context of the inter prediction block.

### C. Cross-Channel Prediction

Cross-channel prediction is to predict between different channels. In the YUV format, the luma channel (Y) is usually coded before the chroma channels (U and V). Thus, it is possible to predict U from Y, and to predict V from Y and U. A traditional method, known as Linear Model (LM), is intended for cross-channel prediction. The key idea of LM is that chroma can be predicted from luma using a linear function, but the coefficients of the function is not transmitted; instead, they are estimated from the context by performing a linear regression. The linear assumption seems over simplified.

Baig and Torresani [7] investigate colorization for image compression. Colorization is to predict chroma from luma, which is an ill-posed problem because one luma value can correspond to multiple chroma values. Accordingly, they propose a tree-structured CNN, which is able to generate multiple predictions (called multiple hypotheses) given one grayscale image as input. When used for compression, the trained CNN is applied at the encoder side, and the branch that produces the best prediction signal is encoded as side information for decoder. They integrate the method into JPEG, without changing the coding of luma, and experimental results show that the proposed method outperforms JPEG for chroma coding.

Li *et al.* [67] propose a cross-channel prediction method analogous to LM. In particular, they design a hybrid neural network consisting of a fully connected part and a convolutional part. The former is used to process the context, including three channels, and the latter is to process the luma channel of the current block. Twofold features are fused to get the final prediction. This method outperforms LM by providing more than 2% BD-rate for chroma coding.

### D. Probability Distribution Prediction

As mentioned before, accurate probability estimation is the key problem in entropy coding. Thus, several works have been done to utilize deep learning for probability distribution prediction to improve the entropy coding efficiency. These works deal with different parts of the information. For example, the intra prediction mode of each block is required to be sent to decoder, and Song *et al.* [106] design a CNN to predict the probability distribution of the intra prediction mode based on the context. Similarly, Pfaff *et al.* [94] predict the probability distribution of the intra prediction mode based on the context, but using a fully connected network. If an encoding/decoding scheme allows multiple transforms and each block can be assigned a transform mode, then Puri *et al.* [96] propose to use a CNN to predict the probability distribution of the transform mode, which is based on the quantized transform coefficients. In a more recent work, Ma *et al.* [82] consider the entropy coding of the quantized transform

coefficients, specifically the DC coefficients. They design a CNN to predict the probability distribution of the DC coefficient of a block, from the context of the block as well as the AC coefficients of the block.

#### E. Transform

Transform is an important tool in the hybrid video coding framework to convert signal (usually residues) into coefficients that are then quantized and coded. At the very beginning, video coding schemes adopt discrete cosine transform (DCT), which is then replaced by integer cosine transform (ICT) in H.264. HEVC also adopts ICT but additionally uses integer sine transform for  $4 \times 4$  luma blocks. Adaptive multiple transforms and secondary transform are also studied. Nonetheless, all these transforms are still very simple.

Inspired by auto-encoder, Liu *et al.* [73] propose a CNN-based method to achieve a DCT-like transform for image coding. The proposed transform consists of a CNN and a fully connected layer, where the CNN is to preprocess the input block and the fully connected layer is to fulfill the *transform*. In their implementation, the fully connected layer is initialized by the transform matrix of DCT, but then is trained together with the CNN. They use a joint rate-distortion cost to train the network, where rate is estimated by the  $l_1$ -norm of the quantized coefficients. They also investigate asymmetric auto-encoders, i.e. the encoding part and decoding part are not symmetric, different from the traditional auto-encoders. Their experimental results show that the trained transform is better than the fixed DCT, and the asymmetric auto-encoders can be useful to achieve a tradeoff between compression efficiency and encoding/decoding time.

#### F. Post- or In-Loop Filtering

Most of the widely used image and video coding schemes are lossy coding ones, i.e. the reconstructed image/video is not exactly the original image/video, for the sake of compression. The loss is usually due to the quantization process shown in Fig. 3. When the quantization step is large, the loss is large too, which may lead to visible artifacts in the reconstructed image/video, such as blocking, blurring, ringing, color shift, and flickering. Filtering is the tool to reduce these artifacts, to improve the quality of the reconstructed image/video, and thus to improve the compression efficiency indirectly. For image, the filtering is also known as post-processing because it does not change the encoding process. For video, the filtering is divided into in-loop and out-of-loop, depending on whether the filtered frame is used as reference for the following frames. In HEVC, two in-loop filters are presented, namely deblocking filter (DF) [91] and sample adaptive offset (SAO) [31].

Post- or in-loop filtering occupies the majority of the related works about deep learning-based image/video coding:

- Earlier works have focused on post-filtering for image coding, especially JPEG. For example, Dong *et al.* [26] propose a 4-layer CNN for compression artifacts reduction, namely ARCNN. ARCNN achieves more than 1dB improvement in PSNR than JPEG on the 5 classical test images when the quality factor (QF) is between 10 and 40. Cavigelli *et al.* [16] use a deeper CNN (12-layer) with hierarchical skip connections and test for higher QF from 40 to 76. Wang *et al.* [125] leverage the prior knowledge of JPEG compression, i.e. quantization of the DCT coefficients of  $8 \times 8$  blocks, and propose a dual-domain (pixel domain and transform domain) based method. They achieve both higher quality and less computing time than ARCNN. Dual-domain processing is

also studied in [38]. Guo and Chao [39] propose a one-to-many network, which is trained by a combination of perceptual loss, naturalness loss, and JPEG loss. Another work about loss function is presented in [32], which suggests the usage of discriminator loss like in GAN. Ororbia *et al.* [92] propose an iterative post-filtering method by using a trained RNN. Recently, several works treat JPEG post-filtering as an image restoration task, like denoising or super-resolution, and propose different networks for a series of image restoration tasks [78], [140], [142], [144].

- Later on, researches are more and more conducted for out-of-loop filtering in video coding, especially HEVC. Dai *et al.* [24] propose a 4-layer CNN for post-filtering of intra frames, where the CNN has variable filter size and residue connection, and named VRCNN. Wang *et al.* [122] use a 10-layer CNN for out-of-loop filtering, where they train a CNN to filter one image and used the trained CNN on the video frames individually. Yang *et al.* [138] propose to train different CNN models for I frames and P frames respectively, and verify the benefit. Jin *et al.* [50] suggest the use of a discriminator loss in addition to the MSE loss. Li *et al.* [62] propose to transmit some side information to decoder to select one model for each frame from a previously trained set of models. In addition, Yang *et al.* [139] propose to utilize the inter-picture correlation during the post-filtering process by inputting multiple neighboring frames into the CNN to enhance one frame. Wang *et al.* [123] also consider the inter-picture correlation, but using a multi-scale convolutional LSTM. While the aforementioned works take only the decoded frames as input to the CNN, He *et al.* [41] propose to input the block partition information together with decoded frame into the CNN, Kang *et al.* [53] also input the block partition information into the CNN and design a multi-scale network, Ma *et al.* [83] input the intra prediction signal and the decoded residual signal into the CNN, and Song *et al.* [107] input the QP plus the decoded frame into the CNN (they also quantize the network parameters to ensure consistency between different computing platforms). A different work is presented in [114], which does not enhance the decoded frames directly; instead, they propose to calculate the compression residues (i.e. the original video minus the decoded video, to be distinguished from prediction residues) at the encoder side, and train an auto-encoder to encode the compression residues and send to the decoder side. Their method is reported to perform well on domain-specific video sequences, e.g. in video game streaming services.
- It is more challenging to integrate CNN-based filter into the coding loop because the filtered frame will serve as reference and will affect the other coding tools. Park and Kim [93] train a 3-layer CNN as an in-loop filter for HEVC. They train two models for two QP ranges: 20–29 and 30–39, respectively, and use one model for each frame according to its QP. The CNN is applied after DF, and SAO is turned off. They also design two cases to apply the CNN-based filter: in the one case, the filter is applied on specified frames based on picture order count (POC); in the other, the filter is tested for each frame and if it improves quality then it is applied, one binary flag for each frame is signaled to decoder in this case. Meng *et al.* [85] use an LSTM as an in-loop filter, which is applied after DF and before SAO in HEVC. The network has decoded frame together with block partition information as its input, and is trained with a combination of MS-SSIM loss and MAE loss. Zhang *et al.* [145] propose a residual highway CNN (RHCNN) for in-loop filtering in HEVC. RHCNN-based filter is applied after SAO. They train different RHCNN models for I, P, and B frames, respectively. They also

divide QPs into several ranges and train a separate model for each range. Dai *et al.* [25] propose a deep CNN called VRCNN-ext for in-loop filtering in HEVC. They design different strategies for I frames and P/B frames: CNN-based filter replaces DF and SAO for I frames, but is applied after DF and before SAO for P/B frames with CTU- and CU-level control. At CTU-level, one binary flag for each CTU is signaled to control the on/off of CNN-based filter; if the flag is off, then at CU-level, a binary classifier is used to decide whether to turn on CNN-based filter for each CU. Jia *et al.* [46] also consider a deep CNN for in-loop filtering in HEVC. The filter is applied after SAO and controlled by frame- and CTU-level flags. If frame-level flag is “off,” then the corresponding CTU-level flags are omitted. In addition, they train multiple CNN models and train a content analysis network that decides one model for each CTU, which saves the bits of CNN model selection.

### G. Down- and Up-Sampling

A trend of the video technology is to increase the resolution at different dimensions, such as spatial resolution (i.e. number of pixels), temporal resolution (i.e. frame rate), and pixel value resolution (i.e. bit-depth). The increasing resolution results in multiplied data volume, which raises a great challenge to video transmission systems. When the bandwidth for transmission is limited (e.g. using 2G or 3G mobile network), a common practice is to decrease video resolution before encoding and to increase video resolution back after decoding. This is known as the down- and up-sampling-based coding strategy. The down- and up-sampling can be performed in the spatial domain, the temporal domain, the pixel value domain, or a combination of these domains. Traditionally, the down- and up-sampling filters are often handcrafted. Recently, it is proposed to train deep networks as down- and up-sampling filters for efficient video coding. There are two categories of related researches.

The first category is focused on training deep networks as up-sampling filters only, while still using handcrafted down-sampling filters. This is inspired by the success of super-resolution, e.g. [27]. For example in [1], a joint spatial and pixel value down-sampling is proposed, where the spatial down-sampling is achieved by a handcrafted low-pass filter and the pixel value down-sampling is achieved by bitwise right shift. At the encoder side, a support vector machine is used to decide whether to perform down-sampling for each frame. At the decoder side, a CNN is trained to up-sample the decoded video to its original resolution. In [69], Li *et al.* only consider spatial down-sampling which is also performed by a handcrafted filter, and train a CNN for up-sampling. But different from [1], they propose a block adaptive resolution coding (BARC) framework. Specifically, for each block inside a frame, they consider two coding modes: down-sampling then coding and direct coding. The encoder can choose a mode for each block and signal the chosen mode to the decoder. In addition, in the down-sampling coding mode, they further design two sub-modes: using a handcrafted simple filter for up-sampling, and using the trained CNN for up-sampling. The sub-mode is also signaled to the decoder. Li *et al.* [69] investigate BARC only for I frames. Later, Lin *et al.* [72] extend the BARC framework for P and B frames and build a complete BARC-based video coding scheme. While the aforementioned works perform down-sampling in the pixel domain, Liu *et al.* [77] propose down-sampling in the residue domain, i.e. they down-sample the inter prediction residues, and they up-sample the residues by a trained CNN with considering the prediction signal. They also follow the BARC framework.

The second category trains not only up-sampling but also down-sampling filters to allow for more flexibility. For example in [47], a compression framework with two CNNs is studied. The first CNN down-samples an image, the down-sampled image is then compressed by an existing image encoder (such as JPEG and BPG), and then decoded, the second CNN up-samples the decoded image. One drawback of this framework is that it cannot be trained end-to-end because the image encoder/decoder is not differentiable. To address this problem, Jiang *et al.* [47] decide to optimize the two CNNs alternatively. Differently, Zhao *et al.* [147] use a virtual codec that is actually a CNN to approximate the functionality of—and thus replace—the encoder/decoder; they also insert a CNN to perform post-processing before the up-sampling CNN; their scheme is fully convolutional and can be trained end-to-end. Moreover, Li *et al.* [68] simply remove the encoder/decoder and keep only the two CNNs during training; considering that the down-sampled image will be compressed, they propose a novel regularization loss for training, which requires the down-sampled image to be not quite different from the ideal low-pass and decimated (which is approximated by a handcrafted filter) image. The regularization loss is verified to be useful when training down- and up-sampling CNNs jointly for image coding.

#### H. Encoding Optimizations

The aforementioned deep tools are intended for increasing the compression efficiency, especially for reducing bitrate while keeping the same PSNR. There are some other deep tools that target different aspects. In this subsection, we review several deep tools for three different objectives: fast encoding, rate control, and region-of-interest (ROI) coding. Since these tools are used only at the encoder side, we call them encoding optimization tools in summary.

1) *Fast Encoding*: Regarding the state-of-the-art video coding standards, H.264 and HEVC, the decoder is computationally simple, but the encoder is much more complex. This is because more and more coding modes are introduced into the video coding standards, and each block can be assigned a different mode. The mode of each block is signaled to the decoder, so the decoder only needs to compute the given mode. But to find the mode for each block, the encoder usually needs to compare the multiple optional modes and select the optimal one, where optimality is claimed in the rate-distortion sense. Therefore, if the encoder performs an exhaustive search, then the compression efficiency is the highest, but the computational complexity may be also very high. Any practical encoder will adopt heuristic algorithms to search for a better mode, where machine learning especially deep learning can help.

Liu *et al.* [79] present a hardware design for HEVC intra encoder, where they adopt a trained CNN to help decide CU partition mode. Specifically in HEVC intra coding, a CTU is split into CUs recursively to form a quadtree structure. Their trained CNN will decide whether to split a  $32 \times 32 / 16 \times 16 / 8 \times 8$  CU or not based on the content inside the CU and the specified QP. Actually, this is a binary decision problem. Xu *et al.* [134] additionally consider HEVC inter encoder, and propose an early-terminated hierarchical CNN and an early-terminated hierarchical LSTM to help decide CU partition mode, for I frames and P/B frames, respectively. Jin *et al.* [49] also consider the CU partition mode decision but for the incoming VVC rather than HEVC, because in VVC a quadtree-bintree (QTBT) structure is designed for CU partition, which is more complex than that in HEVC. They train a CNN to perform 5-way classification for a  $32 \times 32$  CU, where different classes indicate different tree depths. Xu *et al.* [133] investigate

the CU partition mode decision for H.264 to HEVC transcoding. They design a hierarchical LSTM network to predict the CU partition mode from the features extracted from H.264 coded bits.

Song *et al.* [105] study a CNN-based method for fast intra prediction mode decision in HEVC intra encoder. They train a CNN to derive a list of most probable modes for each  $8 \times 8/4 \times 4$  block based on the content and the specified QP, and then choose a mode from the list by the normal rate-distortion optimized process.

2) *Rate Control*: Given a limited transmission bandwidth, video encoder tries to produce bits that do not overflow the bandwidth. This is known as the rate control requirement.

One traditional rate control method is to allocate bits to different blocks according to the  $R-\lambda$  model [61]. In that model, each block has two parameters  $\alpha$  and  $\beta$  that are to be determined. Previously, the parameters are estimated by an empirical formula. In [66], Li *et al.* propose to train a CNN to predict the parameters for each CTU. Experimental results show that the proposed method achieves higher compression efficiency as well as lower rate control error.

Hu *et al.* [43] attempt to leverage reinforcement learning methods for intra-frame rate control. They draw an analogy between the rate control problem and the reinforcement learning problem: the texture complexity of the blocks and bit balance are regarded as the environment state, the quantization parameter is regarded as an action that an agent needs to take, and the negative distortion of the blocks is regarded as an immediate reward. They train a neural network as the agent.

3) *ROI Coding*: ROI refers to the regions of interest in an image. In image compression, it is often required that the content in ROI shall be of high quality and the content not in ROI can be of low quality. Many image coding schemes, such as JPEG and JPEG 2000, support ROI coding. Then, how to identify the ROI is a research problem and has been addressed by using deep learning. Prakash *et al.* [95] propose a CNN-based method to generate a multi-scale ROI (MS-ROI) map to guide the following JPEG coding. They use a trained image classification network on an image, pick the top five classes predicted by the network, and identify the regions that correspond to these classes. Thus, their MS-ROI map indicates salient regions that are related to semantic analysis.

#### IV. CASE STUDY OF DLVC

We now turn to the case study of our developed DLVC, a prototype video codec. Indeed, DLVC was developed as a proposal in response to the joint call for proposals on video compression with capability beyond HEVC. Now the source code of DLVC has been released for future researches<sup>3</sup>. DLVC is crafted upon the JEM software, contains a number of improvements than JEM, and especially has two deep coding tools: CNN-based in-loop filter (CNN-ILF) and CNN-based block adaptive resolution coding (CNN-BARC), both of which are based on trained CNN models. The scheme of DLVC is illustrated in Fig. 5. In this section, we focus on the two deep tools. More technical details about DLVC can be found in the technical report [132].

<sup>3</sup><https://github.com/fvc2018/dlvc>, <http://dlvc.bitahub.com/>.

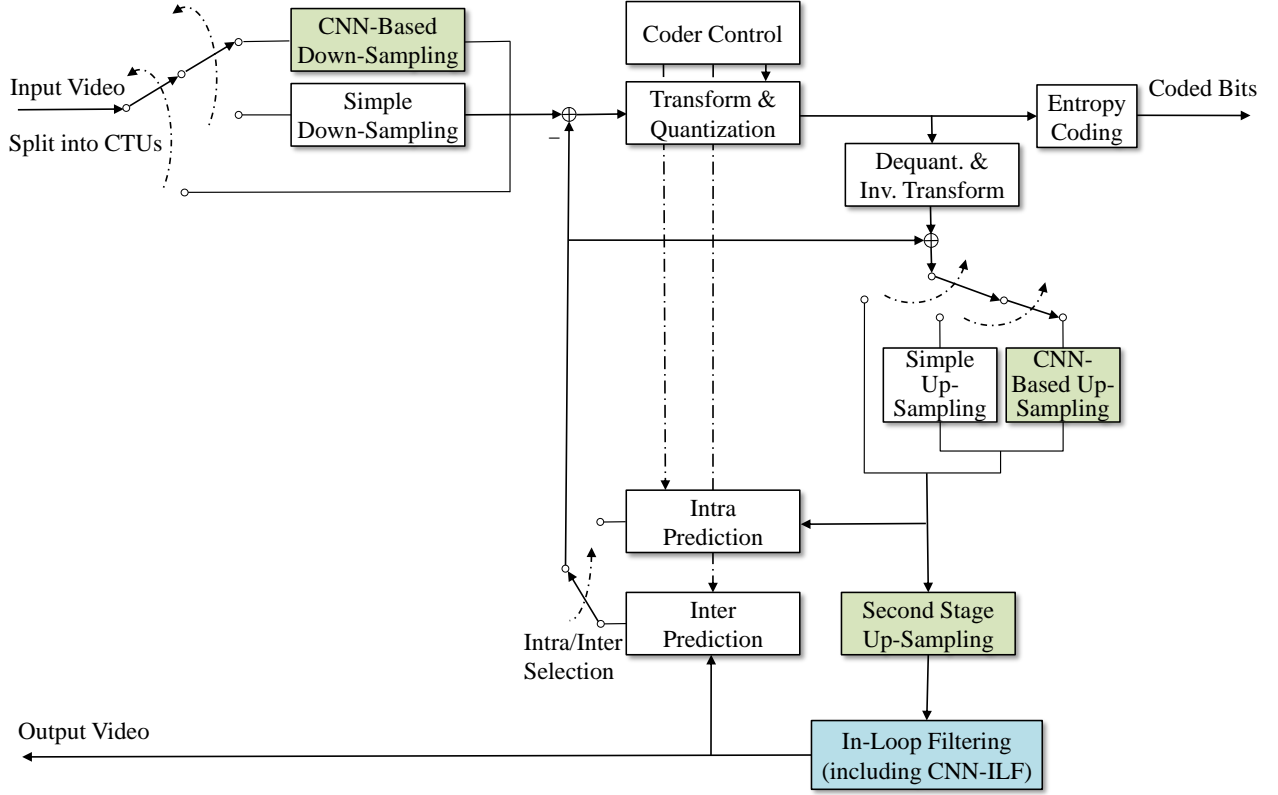


Fig. 5. Illustration of the developed DLVC scheme. Inside the blue block is the proposed CNN-ILF. The green blocks correspond to CNN-BARC.

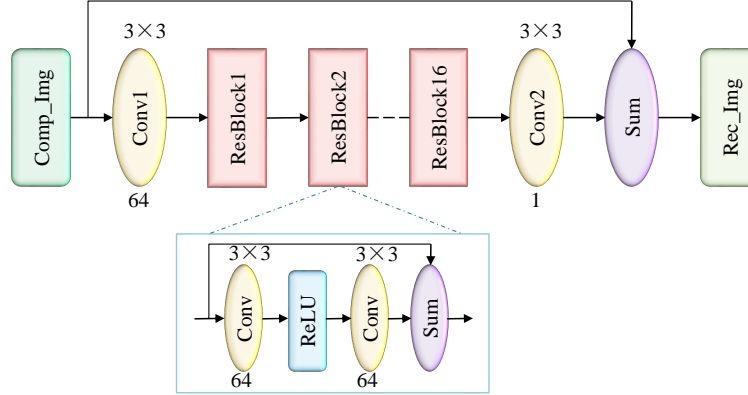


Fig. 6. The network structure of the proposed CNN-ILF. The numbers shown above and beneath each convolutional (Conv) layer indicate the kernel size (e.g.  $3 \times 3$ ) and the number of output channels (e.g. 64), respectively.

#### A. CNN-Based In-Loop Filter

As mentioned in Section III-F, a great number of researches have been conducted on using trained CNN models for post- or in-loop filtering. CNN-ILF represents our efforts at this aspect.

The network structure of our proposed CNN-ILF is illustrated in Fig. 6. Inspired by the SR network in [70],



we design a deep CNN having 16 residual blocks (ResBlocks) and 2 convolutional layers, in total 34 layers. Each ResBlock consists of 2 convolutional layers separated by a ReLU mapping, and a skip connection. The entire network has a global skip connection from the input to the output. These skip connections are crucial to train an efficient network and to accelerate the convergence in training.

To train the network, we have used a set of natural images and compressed each image by the DLVC intra coding (turning off all in-loop filters) at different QPs. For each QP we train a separate model. We only use the luma component for training but the trained models are used for both luma and chroma channels during compression. We divide images into  $70 \times 70$  sub-images and shuffle the sub-images to prepare training data. The loss function is MSE, i.e. the error between the network-output image and the original uncompressed image. We use the stochastic gradient descent algorithm to train the network until convergence.

We apply the trained models in DLVC. The CNN-ILF is applied after deblocking filter and before sample adaptive offset. There are different models corresponding to different QPs, and one model is selected for each frame according to the frame's QP. For each CTU there are two binary flags that control the on/off of the CNN-ILF for luma and chroma respectively. These flags are decided at the encoder side and transmitted to the decoder.

#### *B. CNN-Based Block Adaptive Resolution Coding*

CNN-BARC is a down- and up-sampling-based coding tool that uses trained CNN models as the down- and up-sampling filters. In DLVC, CNN-BARC is applied only for intra frame coding. The down-sampling coding or direct coding mode is decided for each CTU, and the down-sampling coding mode has two sub-modes: using CNNs for down- and up-sampling, and using simple interpolation filters for down- and up-sampling. All the modes and sub-modes are decided by the encoder and signaled to the decoder.

The networks for down- and up-sampling are illustrated in Fig. 7. Specifically, the down-sampling CNN (CNN-DS) has 10 convolutional layers where the first layer has a stride of 2 to achieve  $2 \times$  down-sizing. CNN-DS also embraces residue learning but here the original image is bicubic down-sampled to serve as the skip connection. The up-sampling CNN (CNN-US) is similar to the SR network in [70], and has 16 ResBlocks, 3 convolutional layers, 1 deconvolutional layer, and a global skip connection.

The CNN-DS and CNN-US are trained in four steps. First, we remove the convolutional layers in the CNN-DS, making it a simple bicubic down-sampling, and train the CNN-US to minimize the end-to-end MSE (i.e. the error between original image and down-sampled-up-sampled image). Second, we add back the layers of CNN-DS, fix the parameters of the CNN-US, and train the CNN-DS to minimize the end-to-end MSE. Third, we fine-tune the parameters of CNN-DS and CNN-US simultaneously, using a combination of two losses: the one is end-to-end MSE, and the other is the down-sampled MSE (i.e. the error between bicubic down-sampled image and network down-sampled image), where the latter serves as a regularization term. Fourth, we fix the parameters of the CNN-DS, and compress the down-sampled images by the DLVC intra coding (turning off all in-loop filters) at different QPs. For each QP we train a separate CNN-US model.

There are two mode selection steps regarding CNN-BARC in the DLVC encoder. The first is to decide which down- and up-sampling (sub-)mode, and the second is to decide whether to perform down-sampling. We compare

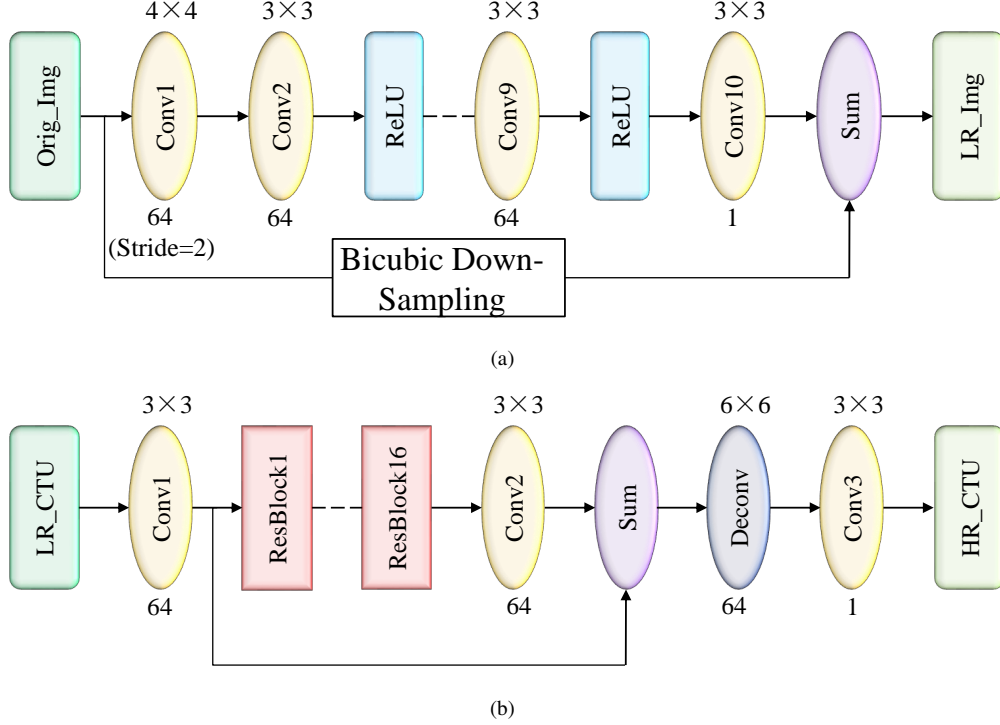


Fig. 7. The network structure of the proposed CNN-BARC, including (a) CNN-DS for down-sampling and (b) CNN-US for up-sampling. Note that the first Conv layer in (a) has a stride of 2 to achieve  $2\times$  down-sampling.

the rate-distortion costs of different modes to make decision. The rate is the number of coded bits, and the distortion is the MSE between original and reconstructed CTUs. For fair comparison, we always calculate the distortion at the original resolution. Last but not the least, after an intra frame is compressed, we perform up-sampling again for the down-sampled-coded CTUs. More details about CNN-BARC can be found in [68], [69].

### C. Compression Performance

We test the compression performance of DLVC on the 10 video sequences that are recommended by JVET. These sequences are divided into Class A and Class B according to spatial resolution: 5 sequences have UHD ( $3840\times 2160$ ) resolution and the other 5 have HD ( $1920\times 1080$ ) resolution. Different coding configurations including all-intra (AI), low-delay (LD), and random-access (RA) are tested. We compare DLVC with the HEVC reference software–HM version 16.16<sup>4</sup> and its variants as well as JEM version 7.0<sup>5</sup>, and use BD-rate [13] to measure the relative compression efficiency.

Table II presents the BD-rate results of DLVC compared to the HEVC anchor. Obviously, DLVC improves the compression efficiency very significantly. Considering the Y channel, DLVC achieves on average 39.6% and 33.0% BD-rate reduction than HEVC, under RA and LD configurations, respectively. The results indicate that when using DLVC to replace HEVC, the bits can be reduced by more than 30% without lowering the reconstruction quality.

<sup>4</sup>HM version 16.16, [https://hevc.hhi.fraunhofer.de/svn/svn\\_HEVCSoftware/tags/HM-16.16/](https://hevc.hhi.fraunhofer.de/svn/svn_HEVCSoftware/tags/HM-16.16/).

<sup>5</sup>JEM version 7.0, [https://jvet.hhi.fraunhofer.de/svn/svn\\_HMJEMSoftware/tags/HM-16.6-JEM-7.0/](https://jvet.hhi.fraunhofer.de/svn/svn_HMJEMSoftware/tags/HM-16.6-JEM-7.0/).

TABLE II  
BD-RATE RESULTS OF DLVC COMPARED TO HM16.16

	Random-Access			Low-Delay		
	Y	U	V	Y	U	V
FoodMarket	−38.42%	−45.57%	−50.10%	-	-	-
CatRobot	−46.87%	−59.26%	−53.80%	-	-	-
DaylightRoad	−47.34%	−60.65%	−46.19%	-	-	-
ParkRunning	−38.43%	−29.65%	−31.58%	-	-	-
CampfireParty	−41.20%	−38.07%	−54.53%	-	-	-
Avg. Class-A	−42.45%	−46.64%	−47.24%	-	-	-
BQTerrace	−37.53%	−50.06%	−56.97%	−35.56%	−52.79%	−59.11%
RitualDance	−33.30%	−42.15%	−45.71%	−28.23%	−41.19%	−43.19%
MarketPlace	−35.34%	−51.78%	−49.99%	−26.85%	−46.65%	−47.83%
BasketballDrive	−37.09%	−52.86%	−51.20%	−32.98%	−50.91%	−52.69%
Cactus	−40.80%	−50.92%	−46.82%	−41.37%	−56.66%	−54.68%
Avg. Class-B	−36.81%	−49.55%	−50.14%	−33.00%	−49.64%	−51.50%
<b>Avg. All</b>	<b>−39.63%</b>	<b>−48.10%</b>	<b>−48.69%</b>	<b>−33.00%</b>	<b>−49.64%</b>	<b>−51.50%</b>

Table III presents the BD-rate results of DLVC (document number J0032) compared to the JEM anchor. For comparison, we also include the BD-rate results of the other proposals in response to the joint call for proposals. Considering the Y channel, DLVC achieves on average 10.1% and 11.8% BD-rate reduction than JEM, under RA and LD configurations, respectively. DLVC is among the best proposals from the BD-rate perspective.

Table IV verifies the effectiveness of the proposed CNN-ILF. Specifically, we use a variant of HM that adds the QTBT structure, which outperforms the vanilla HM, as the anchor. We integrate the CNN-ILF into the anchor and turn on/off the CNN-ILF for comparison. As shown, CNN-ILF achieves significant BD-rate reduction: on average 5.5%, 5.2%, 6.4% for the Y channel under RA, LD, AI configurations, respectively.

Table V verifies the effectiveness of the proposed CNN-BARC. We use another variant of HM that adds the quadtree-bintree-triple-tree (QTBTTT) structure, which further outperforms the HM plus QTBT, as the anchor. We integrate the CNN-BARC into the anchor and turn on/off the CNN-BARC for comparison. As shown, CNN-BARC achieves significant BD-rate reduction under AI configuration: on average 5.4% for the Y channel. The BD-rate under RA and LD configurations is less significant, because CNN-BARC is applied on intra frames only.

## V. PERSPECTIVES AND CONCLUSIONS

### A. Open Problems

- *Deep Schemes or Deep Tools.* Shall we be ambitious to expect deep scheme to be the future of video coding, or shall we be satisfied with deep tools within traditional non-deep schemes? In other words, can the non-deep schemes be completely replaced by deep schemes? As for now the answer to this question is probably “no” because deep schemes in general do not outperform non-deep schemes for video coding. But as research goes on, the answer may become “yes” via two ways: first, deep schemes may be improved so much that they

TABLE III

BD-RATE RESULTS OF ALL THE PROPOSALS IN RESPONSE TO THE JOINT CALL FOR PROPOSALS COMPARED TO JEM7.0

Proposal	Organizations	Random-Access			Low-Delay		
		Y	U	V	Y	U	V
J0011	DJI, Peking Univ.	−1.57%	−0.71%	−1.72%	−3.30%	−0.67%	−4.26%
J0012	Ericsson, Nokia	−0.90%	−1.14%	−1.15%	−0.17%	−0.48%	−1.93%
J0013	ETRI, Sejong Univ.	0.64%	−0.39%	−0.89%	0.85%	−1.27%	−2.50%
J0014	Fraunhofer HHI	−7.55%	−6.94%	−5.96%	−7.22%	−7.62%	−5.71%
J0015	InterDigital, Dolby	−3.98%	−3.28%	−3.16%	−3.64%	−2.55%	−4.48%
J0016	KDDI	−0.57%	−0.52%	−1.30%	−0.17%	0.38%	−0.40%
J0017	LG Electronics	−2.52%	−5.29%	−6.19%	−2.09%	−2.47%	−3.70%
J0018	MediaTek	−16.60%	−6.75%	−10.43%	−9.41%	−1.92%	−3.35%
		−14.40%	−5.13%	−8.82%	−7.20%	0.23%	−0.96%
J0020	Panasonic	−2.28%	−3.44%	−3.88%	−2.02%	−1.02%	−2.36%
		−0.06%	0.91%	0.52%	−0.09%	2.87%	1.21%
J0021	Qualcomm, Technicolor	−15.53%	−3.66%	−5.97%	−12.65%	−17.40%	−19.95%
		−10.26%	0.05%	−1.65%	−9.87%	−11.73%	−14.88%
J0022		−13.60%	−3.80%	−5.63%	−12.69%	−16.65%	−18.75%
J0023	RWTH Aachen Univ.	−0.79%	−1.52%	−1.52%	−0.84%	−0.58%	−0.80%
J0024	Samsung, Huawei, GoPro, HiSilicon	−6.01%	10.34%	8.53%	−0.38%	14.73%	15.84%
		−4.24%	10.71%	9.23%	—	—	—
J0026	Sharp, Foxconn	−6.15%	−5.68%	−5.68%	−5.57%	9.00%	8.69%
J0027	NHK, Sharp	−2.14%	−5.55%	−5.61%	−2.23%	−1.97%	−3.83%
J0028	Sony	−2.41%	−4.85%	−5.14%	−2.25%	−6.74%	−7.34%
J0029	Tencent	−4.70%	−8.34%	−8.91%	−4.47%	−9.47%	−10.82%
J0031	Univ. Bristol	−4.54%	20.19%	18.68%	−0.52%	5.74%	5.81%
<b>J0032</b>	<b>USTC and others</b>	<b>−10.11%</b>	<b>−9.59%</b>	<b>−9.97%</b>	<b>−11.83%</b>	<b>−13.22%</b>	<b>−16.49%</b>

More details can be obtained at

[http://phenix.it-sudparis.eu/jvet/doc\\_end\\_user/current\\_meeting.php?id\\_meeting=174&search\\_id\\_group=1&search\\_sub\\_group=1](http://phenix.it-sudparis.eu/jvet/doc_end_user/current_meeting.php?id_meeting=174&search_id_group=1&search_sub_group=1).

TABLE IV

BD-RATE RESULTS OF CNN-ILF ON TOP OF HM16.16 + QTBT

	Random-Access			Low-Delay			All-Intra		
	Y	U	V	Y	U	V	Y	U	V
Avg. Class-A	−5.1%	−3.8%	−4.5%	-	-	-	−5.7%	−5.1%	−5.9%
Avg. Class-B	−5.9%	−5.8%	−5.8%	−5.2%	−7.2%	−8.4%	−7.1%	−6.5%	−7.8%
<b>Avg. All</b>	<b>−5.5%</b>	<b>−4.8%</b>	<b>−5.1%</b>	<b>−5.2%</b>	<b>−7.2%</b>	<b>−8.4%</b>	<b>−6.4%</b>	<b>−5.8%</b>	<b>−6.8%</b>

clearly beat non-deep schemes; second, the coding tools in a traditional coding scheme (e.g. HEVC) may be all replaced by corresponding deep tools, leading to a “deepened” coding scheme that is better than before. The second way may be more practical according to our subjective feeling.

- *Compression Efficiency versus Computational Complexity.* Compared the existing deep tools with their coun-

TABLE V  
BD-RATE RESULTS OF CNN-BARC ON TOP OF HM16.16 + QTBT

	Random-Access			Low-Delay			All-Intra		
	Y	U	V	Y	U	V	Y	U	V
Avg. Class-A	−1.8%	0.6%	0.4%	-	-	-	−7.2%	2.0%	2.0%
Avg. Class-B	−1.0%	0.1%	−0.1%	−0.3%	−0.3%	−0.0%	−3.6%	0.5%	0.4%
<b>Avg. All</b>	<b>−1.4%</b>	<b>0.3%</b>	<b>0.2%</b>	<b>−0.3%</b>	<b>−0.3%</b>	<b>−0.0%</b>	<b>−5.4%</b>	<b>1.3%</b>	<b>1.2%</b>

terparts in the traditional non-deep schemes, one may easily notice that the computational complexity of the former is much higher than the latter. High complexity is indeed a general issue of deep learning, and a critical issue that hinders the adoption of deep networks in scenarios of limited computing resource, e.g. mobile phones. This general issue is now addressed at two aspects: first, to develop novel, efficient, compact deep networks that maintain the high performance (i.e. compression efficiency for video coding) but require much less computations; second, to advocate the adoption of hardware that is specifically designed for deep networks.

- *Optimization for Perceptual Naturalness or Semantic Quality.* Coding schemes designed for natural video are usually serving for human viewing, e.g. television, movie, micro-video. It is natural for these schemes that the quality of reconstructed video shall be evaluated based on human perception. Nonetheless, for traditional non-deep coding schemes, the most widely adopted quality metric is still PSNR, which corresponds to human perception poorly. For deep schemes or deep tools, a few works have been done to optimize them for perceptual naturalness, e.g. using discriminator loss. Moreover, there are coding schemes that serve for automatic semantic analysis instead of human viewing, such as surveillance video coding. For these schemes, the quality metric shall be semantic quality [74], which remains largely unexplored. As a special note, we find that there is a tradeoff between signal fidelity, perceptual naturalness, and semantic quality [75], which implies that the optimization target shall be aligned with the actual requirement.
- *Speciality and Universality.* To one extreme, can one coding scheme be simply the best for any kind of video? The answer is “no” due to the no-free-lunch theorem, which is claimed in the machine learning literature [130] and also applies for compression. To another extreme, can we have a special coding scheme for each video? Not to mention the practical difficulty, such a coding “strategy” is useless because it is no more than assigning an identifier to each video. In between the two extremes are practical and useful coding schemes. That says, coding schemes shall have both speciality and universality to some extent. For deep schemes and deep tools, it implies that the training data shall be carefully selected to reflect the interesting data domain. Researches at this aspect are expected.
- *Federated Design of Multiple Deep Tools.* Currently, most of deep tools have been designed individually, but once they are applied jointly, they may not collaborate well, or may even conflict with each other. The underlying reason is that multiple coding tools are indeed dependent. For example, different prediction tools generate different predictions and lead to variety of residual signal, so transform tools dealing with residual

signal perform differently. Ideally, multiple deep tools shall be designed in a federated manner. However, this can be difficult because the dependency among tools is complicated.

### B. Future Work

In the predictable future, the requirement about video coding technology is still increasing. For entertainment, virtual reality and augmented reality applications are calling for techniques to tackle with new data, such as depth map, point cloud, 3D surface, and so on. For surveillance, the need of intelligent analysis pushes the upgrade of video resolution. For scientific observation, more and more observing instruments are directly connected to a video recorder and generate massive video data. All these requirements drive the development of video coding to achieve higher compression efficiency, lower computational complexity, and smarter integration into video analytical systems. We believe deep learning-based video coding techniques are promising for these challenging objectives. Especially, we expect a holistic framework based on deep networks and integrating image/video acquisition, coding, processing, analysis, and understanding, which indeed mimics human vision system.

### ACKNOWLEDGMENT

This work is supported by the Natural Science Foundation of China under Grant 61772483. Authors would like to thank the following colleagues and collaborators: Jizheng Xu, Bin Li, Zhibo Chen, Li Li, Fangdong Chen, Yuanying Dai, Changsheng Gao, Lei Guo, Shuai Huo, Ye Li, Kang Liu, Changyue Ma, Haichuan Ma, Rui Song, Yefei Wang, Ning Yan, Kun Yang, Qingyu Zhang, Zhenxin Zhang, and Haitao Yang.

### REFERENCES

- [1] M. Afonso, F. Zhang, and D. R. Bull, "Video compression based on spatio-temporal resolution adaptation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 1, pp. 275–280, 2019.
- [2] E. Agustsson, F. Mentzer, M. Tschannen, L. Cavigelli, R. Timofte, L. Benini, and L. V. Gool, "Soft-to-hard vector quantization for end-to-end learning compressible representations," in *NIPS*, 2017, pp. 1141–1151.
- [3] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. V. Gool, "Extreme learned image compression with GANs," in *CVPR Workshops*, 2018, pp. 2587–2590.
- [4] E. Ahanonu, "Lossless image compression using reversible integer wavelet transforms and convolutional neural networks," Master's thesis, University of Arizona, 2018.
- [5] M. Akbari, J. Liang, and J. Han, "DSSLIC: Deep semantic segmentation-based layered image compression," in *ICASSP*, 2019, pp. 2042–2046.
- [6] M. H. Baig, V. Koltun, and L. Torresani, "Learning to inpaint for image compression," in *NIPS*, 2017, pp. 1246–1255.
- [7] M. H. Baig and L. Torresani, "Multiple hypothesis colorization and its application to image compression," *Computer Vision and Image Understanding*, vol. 164, pp. 111–123, 2017.
- [8] J. Ballé, "Efficient nonlinear transforms for lossy image compression," in *PCS*, 2018, pp. 248–252.
- [9] J. Ballé, V. Laparra, and E. P. Simoncelli, "End-to-end optimization of nonlinear transform codes for perceptual quality," in *PCS*. IEEE, 2016, pp. 1–5.
- [10] —, "End-to-end optimized image compression," *arXiv preprint arXiv:1611.01704*, 2016.
- [11] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," *arXiv preprint arXiv:1802.01436*, 2018.
- [12] Y. Bengio and S. Bengio, "Modeling high-dimensional discrete data with multi-layer neural networks," in *NIPS*, 2000, pp. 400–406.
- [13] G. Bjontegaard, "Calculation of average PSNR differences between RD-curves," VCEG, Tech. Rep. VCEG-M33, 2001.

- [14] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in *CVPR*, 2018, pp. 6228–6237.
- [15] C. Cai, L. Chen, X. Zhang, and Z. Gao, "Efficient variable rate image compression with multi-scale decomposition network," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2018.2880492, 2018.
- [16] L. Cavigelli, P. Hager, and L. Benini, "CAS-CNN: A deep convolutional neural network for image compression artifact suppression," in *IJCNN*. IEEE, 2017, pp. 752–759.
- [17] T. Chen, H. Liu, Q. Shen, T. Yue, X. Cao, and Z. Ma, "DeepCoder: A deep neural network based video compression," in *VCIP*. IEEE, 2017, pp. 1–4.
- [18] X. Chen, N. Mishra, M. Rohaninejad, and P. Abbeel, "PixelSNAIL: An improved autoregressive generative model," in *ICML*, 2018, pp. 863–871.
- [19] Z. Chen and T. He, "Learning based facial image compression with semantic fidelity metric," *Neurocomputing*, vol. 338, pp. 16–25, 2019.
- [20] Z. Chen, T. He, X. Jin, and F. Wu, "Learning for video compression," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2019.2892608, 2019.
- [21] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, "Deep convolutional autoencoder-based lossy image compression," in *PCS*. IEEE, 2018, pp. 253–257.
- [22] M. Covell, N. Johnston, D. Minnen, S. J. Hwang, J. Shor, S. Singh, D. Vincent, and G. Toderici, "Target-quality image compression with recurrent, convolutional neural networks," *arXiv preprint arXiv:1705.06687*, 2017.
- [23] W. Cui, T. Zhang, S. Zhang, F. Jiang, W. Zuo, Z. Wan, and D. Zhao, "Convolutional neural networks based intra prediction for HEVC," in *DCC*. IEEE, 2017, p. 436.
- [24] Y. Dai, D. Liu, and F. Wu, "A convolutional neural network approach for post-processing in HEVC intra coding," in *MMM*. Springer, 2017, pp. 28–39.
- [25] Y. Dai, D. Liu, Z.-J. Zha, and F. Wu, "A CNN-based in-loop filter with CU classification for HEVC," in *VCIP*, 2018, pp. 1–4.
- [26] C. Dong, Y. Deng, C. C. Loy, and X. Tang, "Compression artifacts reduction by a deep convolutional network," in *ICCV*, 2015, pp. 576–584.
- [27] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *ECCV*. Springer, 2014, pp. 184–199.
- [28] R. D. Dony and S. Haykin, "Neural network approaches to image compression," *Proceedings of the IEEE*, vol. 83, no. 2, pp. 288–303, 1995.
- [29] T. Dumas, A. Roumy, and C. Guillemot, "Image compression with stochastic winner-take-all auto-encoder," in *ICASSP*. IEEE, 2017, pp. 1512–1516.
- [30] —, "Autoencoder based image compression: can the learning be quantization independent?" in *ICASSP*. IEEE, 2018, pp. 1188–1192.
- [31] C.-M. Fu, E. Alshina, A. Alshin, Y.-W. Huang, C.-Y. Chen, C.-Y. Tsai, C.-W. Hsu, S.-M. Lei, J.-H. Park, and W.-J. Han, "Sample adaptive offset in the HEVC standard," *IEEE Transactions on Circuits and Systems for Video technology*, vol. 22, no. 12, pp. 1755–1764, 2012.
- [32] L. Galteri, L. Seidenari, M. Bertini, and A. Del Bimbo, "Deep generative adversarial compression artifact removal," in *ICCV*, 2017, pp. 4826–4835.
- [33] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014, pp. 580–587.
- [34] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014, pp. 2672–2680.
- [35] K. Gregor, F. Besse, D. J. Rezende, I. Danihelka, and D. Wierstra, "Towards conceptual compression," in *NIPS*, 2016, pp. 3549–3557.
- [36] K. Gregor, I. Danihelka, A. Graves, D. Rezende, and D. Wierstra, "DRAW: A recurrent neural network for image generation," in *ICML*, 2015, pp. 1462–1471.
- [37] K. Gregor and Y. LeCun, "Learning representations by maximizing compression," *arXiv preprint arXiv:1108.1169*, 2011.
- [38] J. Guo and H. Chao, "Building dual-domain representations for compression artifacts reduction," in *ECCV*. Springer, 2016, pp. 628–644.
- [39] —, "One-to-many network for visually pleasing compression artifacts reduction," in *CVPR*, 2017, pp. 3038–3047.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [41] X. He, Q. Hu, X. Zhang, C. Zhang, W. Lin, and X. Han, "Enhancing HEVC compressed videos with a partition-masked convolutional neural network," in *ICIP*. IEEE, 2018, pp. 216–220.
- [42] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

- [43] J.-H. Hu, W.-H. Peng, and C.-H. Chung, "Reinforcement learning for HEVC/H.265 intra-frame rate control," in *ISCAS*. IEEE, 2018, pp. 1–5.
- [44] Y. Hu, W. Yang, M. Li, and J. Liu, "Progressive spatial recurrent neural network for intra prediction," *arXiv preprint arXiv:1807.02232*, 2018.
- [45] S. Huo, D. Liu, F. Wu, and H. Li, "Convolutional neural network-based motion compensation refinement for video coding," in *ISCAS*, 2018, pp. 1–4.
- [46] C. Jia, S. Wang, X. Zhang, S. Wang, J. Liu, S. Pu, and S. Ma, "Content-aware convolutional neural network for in-loop filtering in high efficiency video coding," *IEEE Transactions on Image Processing*, DOI: 10.1109/TIP.2019.2896489, 2019.
- [47] F. Jiang, W. Tao, S. Liu, J. Ren, X. Guo, and D. Zhao, "An end-to-end compression framework based on convolutional neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 10, pp. 3007–3018, 2018.
- [48] J. Jiang, "Image compression with neural networks—A survey," *Signal Processing: Image Communication*, vol. 14, no. 9, pp. 737–760, 1999.
- [49] Z. Jin, P. An, L. Shen, and C. Yang, "CNN oriented fast QTBT partition algorithm for JVET intra coding," in *VCIP*. IEEE, 2017, pp. 1–4.
- [50] Z. Jin, P. An, C. Yang, and L. Shen, "Quality enhancement for intra frame coding via CNNs: An adversarial approach," in *ICASSP*. IEEE, 2018, pp. 1368–1372.
- [51] N. Johnston, D. Vincent, D. Minnen, M. Covell, S. Singh, T. Chinen, S. Jin Hwang, J. Shor, and G. Toderici, "Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks," in *CVPR*, 2018, pp. 4385–4393.
- [52] N. Kalchbrenner, A. van den Oord, K. Simonyan, I. Danihelka, O. Vinyals, A. Graves, and K. Kavukcuoglu, "Video pixel networks," in *ICML*, 2017, pp. 1771–1779.
- [53] J. Kang, S. Kim, and K. M. Lee, "Multi-modal/multi-scale convolutional neural network based in-loop filter design for next generation video codec," in *ICIP*, 2017, pp. 26–30.
- [54] S. Kim, J. S. Park, C. G. Bampis, J. Lee, M. K. Markey, A. G. Dimakis, and A. C. Bovik, "Adversarial video compression guided by soft edge detection," *arXiv preprint arXiv:1811.10673*, 2018.
- [55] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [56] A. Kolesnikov and C. H. Lampert, "Latent variable PixelCNNs for natural image modeling," *arXiv preprint arXiv:1612.08185*, 2016.
- [57] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.
- [58] H. Larochelle and I. Murray, "The neural autoregressive distribution estimator," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 29–37.
- [59] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- [60] J. Lee, S. Cho, and S.-K. Beack, "Context-adaptive entropy model for end-to-end optimized image compression," *arXiv preprint arXiv:1809.10452*, 2018.
- [61] B. Li, H. Li, L. Li, and J. Zhang, " $\lambda$  domain rate control algorithm for high efficiency video coding," *IEEE Transactions on Image Processing*, vol. 23, no. 9, pp. 3841–3854, 2014.
- [62] C. Li, L. Song, R. Xie, and W. Zhang, "CNN based post-processing to improve HEVC," in *ICIP*. IEEE, 2017, pp. 4577–4580.
- [63] J. Li, B. Li, J. Xu, R. Xiong, and W. Gao, "Fully connected network-based intra prediction for image coding," *IEEE Transactions on Image Processing*, vol. 27, no. 7, pp. 3236–3247, 2018.
- [64] M. Li, S. Gu, D. Zhang, and W. Zuo, "Enlarging context with low cost: Efficient arithmetic coding with trimmed convolution," *arXiv preprint arXiv:1801.04662*, 2018.
- [65] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, "Learning convolutional networks for content-weighted image compression," in *CVPR*, 2018, pp. 673–681.
- [66] Y. Li, B. Li, D. Liu, and Z. Chen, "A convolutional neural network-based approach to rate control in HEVC intra coding," in *VCIP*. IEEE, 2017, pp. 1–4.
- [67] Y. Li, L. Li, Z. Li, J. Yang, N. Xu, D. Liu, and H. Li, "A hybrid neural network for chroma intra prediction," in *ICIP*, 2018, pp. 1797–1801.
- [68] Y. Li, D. Liu, H. Li, L. Li, Z. Li, and F. Wu, "Learning a convolutional neural network for image compact-resolution," *IEEE Transactions on Image Processing*, vol. 28, no. 3, pp. 1092–1107, 2019.



- [69] Y. Li, D. Liu, H. Li, L. Li, F. Wu, H. Zhang, and H. Yang, "Convolutional neural network-based block up-sampling for intra frame coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 9, pp. 2316–2330, 2018.
- [70] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *CVPR Workshops*, 2017, pp. 136–144.
- [71] J. Lin, D. Liu, H. Li, and F. Wu, "Generative adversarial network-based frame extrapolation for video coding," in *VCIP*, 2018, pp. 1–4.
- [72] J. Lin, D. Liu, H. Yang, H. Li, and F. Wu, "Convolutional neural network-based block up-sampling for HEVC," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2018.2884203, 2018.
- [73] D. Liu, H. Ma, Z. Xiong, and F. Wu, "CNN-based DCT-like transform for image compression," in *MMM*. Springer, 2018, pp. 61–72.
- [74] D. Liu, D. Wang, and H. Li, "Recognizable or not: Towards image semantic quality assessment for compression," *Sensing and Imaging*, vol. 18, no. 1, pp. 1–20, 2017.
- [75] D. Liu, H. Zhang, and Z. Xiong, "On the classification-distortion-perception tradeoff," *arXiv preprint arXiv:1904.08816*, 2019.
- [76] J. Liu, S. Xia, W. Yang, M. Li, and D. Liu, "One-for-all: Grouped variation network based fractional interpolation in video coding," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2140–2151, 2019.
- [77] K. Liu, D. Liu, H. Li, and F. Wu, "Convolutional neural network-based residue super-resolution for video coding," in *VCIP*, 2018, pp. 1–4.
- [78] P. Liu, H. Zhang, K. Zhang, L. Lin, and W. Zuo, "Multi-level wavelet-CNN for image restoration," in *CVPR Workshops*, 2018, pp. 773–782.
- [79] Z. Liu, X. Yu, Y. Gao, S. Chen, X. Ji, and D. Wang, "CU partition mode decision for HEVC hardwired intra encoder using convolution neural network," *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 5088–5103, 2016.
- [80] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, "DVC: An end-to-end deep video compression framework," in *CVPR*, 2019.
- [81] S. Luo, Y. Yang, Y. Yin, C. Shen, Y. Zhao, and M. Song, "DeepSIC: Deep semantic image compression," in *International Conference on Neural Information Processing*. Springer, 2018, pp. 96–106.
- [82] C. Ma, D. Liu, X. Peng, and F. Wu, "Convolutional neural network-based arithmetic coding of DC coefficients for HEVC intra coding," in *ICIP*, 2018, pp. 1772–1776.
- [83] L. Ma, Y. Tian, and T. Huang, "Residual-based video restoration for HEVC intra coding," in *BigMM*. IEEE, 2018, pp. 1–7.
- [84] S. Ma, X. Zhang, C. Jia, Z. Zhao, S. Wang, and S. Wang, "Image and video compression with neural networks: A review," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2019.2910119, 2019.
- [85] X. Meng, C. Chen, S. Zhu, and B. Zeng, "A new HEVC in-loop filter based on multi-channel long-short-term dependency residual networks," in *DCC*. IEEE, 2018, pp. 187–196.
- [86] F. Mentzer, E. Agustsson, M. Tschanen, R. Timofte, and L. V. Gool, "Conditional probability models for deep image compression," in *CVPR*, 2018, pp. 4394–4402.
- [87] —, "Practical full resolution learned lossless image compression," *arXiv preprint arXiv:1811.12817*, 2018.
- [88] D. Minnen, J. Ballé, and G. Toderici, "Joint autoregressive and hierarchical priors for learned image compression," in *NIPS*, 2018, pp. 10 794–10 803.
- [89] D. Minnen, G. Toderici, M. Covell, T. Chinen, N. Johnston, J. Shor, S. J. Hwang, D. Vincent, and S. Singh, "Spatially adaptive image compression using a tiled deep network," in *ICIP*. IEEE, 2017, pp. 2796–2800.
- [90] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010, pp. 807–814.
- [91] A. Norkin, G. Bjontegaard, A. Fuldseth, M. Narroschke, M. Ikeda, K. Andersson, M. Zhou, and G. van der Auwera, "HEVC deblocking filter," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1746–1754, 2012.
- [92] A. G. Ororbia, A. Mali, J. Wu, S. O'Connell, D. Miller, and C. L. Giles, "Learned neural iterative decoding for lossy image compression systems," *arXiv preprint arXiv:1803.05863*, 2018.
- [93] W.-S. Park and M. Kim, "CNN-based in-loop filtering for coding efficiency improvement," in *IEEE Image, Video, and Multidimensional Signal Processing Workshop*. IEEE, 2016, pp. 1–5.
- [94] J. Pfaff, P. Helle, D. Maniry, S. Kaltenstadler, W. Samek, H. Schwarz, D. Marpe, and T. Wiegand, "Neural network based intra prediction for video coding," in *Applications of Digital Image Processing XLI*, vol. 10752. International Society for Optics and Photonics, 2018, p. 1075213.
- [95] A. Prakash, N. Moran, S. Garber, A. DiLillo, and J. Storer, "Semantic perceptual image compression using deep convolution networks," in *DCC*. IEEE, 2017, pp. 250–259.

- [96] S. Puri, S. Lasserre, and P. Le Callet, "CNN-based transform index prediction in multiple transforms framework to assist entropy coding," in *EUSIPCO*. IEEE, 2017, pp. 798–802.
- [97] O. Rippel and L. Bourdev, "Real-time adaptive image compression," in *ICML*, 2017, pp. 2922–2930.
- [98] O. Rippel, S. Nair, C. Lew, S. Branson, A. G. Anderson, and L. Bourdev, "Learned video compression," *arXiv preprint arXiv:1811.06981*, 2018.
- [99] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma, "PixelCNN++: Improving the PixelCNN with discretized logistic mixture likelihood and other modifications," *arXiv preprint arXiv:1701.05517*, 2017.
- [100] S. Santurkar, D. Budden, and N. Shavit, "Generative compression," in *PCS*. IEEE, 2018, pp. 258–262.
- [101] I. Schioppa, Y. Liu, and A. Munteanu, "CNN-based prediction for lossless coding of photographic images," in *PCS*. IEEE, 2018, pp. 16–20.
- [102] C. E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [103] A. Skodras, C. Christopoulos, and T. Ebrahimi, "The JPEG 2000 still image compression standard," *IEEE Signal Processing Magazine*, vol. 18, no. 5, pp. 36–58, 2001.
- [104] J. Snell, K. Ridgeway, R. Liao, B. D. Roads, M. C. Mozer, and R. S. Zemel, "Learning to generate images with perceptual similarity metrics," in *ICIP*. IEEE, 2017, pp. 4277–4281.
- [105] N. Song, Z. Liu, X. Ji, and D. Wang, "CNN oriented fast PU mode decision for HEVC hardwired intra encoder," in *GlobalSIP*. IEEE, 2017, pp. 239–243.
- [106] R. Song, D. Liu, H. Li, and F. Wu, "Neural network-based arithmetic coding of intra prediction modes in HEVC," in *VCIP*, 2017, pp. 1–4.
- [107] X. Song, J. Yao, L. Zhou, L. Wang, X. Wu, D. Xie, and S. Pu, "A practical convolutional neural network as loop filter for intra frame," in *ICIP*. IEEE, 2018, pp. 1133–1137.
- [108] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (HEVC) standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [109] L. Theis and M. Bethge, "Generative image modeling using spatial LSTMs," in *NIPS*, 2015, pp. 1927–1935.
- [110] L. Theis, W. Shi, A. Cunningham, and F. Huszár, "Lossy image compression with compressive autoencoders," *arXiv preprint arXiv:1703.00395*, 2017.
- [111] G. Toderici, S. M. O'Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar, "Variable rate image compression with recurrent neural networks," *arXiv preprint arXiv:1511.06085*, 2015.
- [112] G. Toderici, D. Vincent, N. Johnston, S. J. Hwang, D. Minnen, J. Shor, and M. Covell, "Full resolution image compression with recurrent neural networks," in *CVPR*, 2017, pp. 5306–5314.
- [113] R. Torfason, F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. V. Gool, "Towards image understanding from deep compression without decoding," *arXiv preprint arXiv:1803.06131*, 2018.
- [114] Y.-H. Tsai, M.-Y. Liu, D. Sun, M.-H. Yang, and J. Kautz, "Learning binary residual representations for domain-specific video streaming," in *AAAI*, 2018, pp. 7363–7370.
- [115] P. Tudor, "MPEG-2 video compression," *Electronics & Communication Engineering Journal*, vol. 7, no. 6, pp. 257–264, 1995.
- [116] B. Uria, I. Murray, and H. Larochelle, "RNADE: The real-valued neural autoregressive density-estimator," in *NIPS*, 2013, pp. 2175–2183.
- [117] —, "A deep and tractable density estimator," in *ICML*, 2014, pp. 467–475.
- [118] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," in *ICML*, 2016, pp. 1747–1756.
- [119] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu, "Conditional image generation with PixelCNN decoders," in *NIPS*, 2016, pp. 4790–4798.
- [120] A. van den Oord and B. Schrauwen, "Factoring variations in natural images with deep Gaussian mixture models," in *NIPS*, 2014, pp. 3518–3526.
- [121] G. K. Wallace, "The JPEG still picture compression standard," *IEEE Transactions on Consumer Electronics*, vol. 38, no. 1, pp. xviii–xxxiv, 1992.
- [122] T. Wang, M. Chen, and H. Chao, "A novel deep learning-based method of improving coding efficiency from the decoder-end for HEVC," in *DCC*. IEEE, 2017, pp. 410–419.
- [123] T. Wang, W. Xiao, M. Chen, and H. Chao, "The multi-scale deep decoder for the standard HEVC bitstreams," in *DCC*. IEEE, 2018, pp. 197–206.
- [124] Y. Wang, X. Fan, C. Jia, D. Zhao, and W. Gao, "Neural network based inter prediction for HEVC," in *ICME*. IEEE, 2018, pp. 1–6.

- [125] Z. Wang, D. Liu, S. Chang, Q. Ling, Y. Yang, and T. S. Huang, "D3: Deep dual-domain based fast restoration of JPEG-compressed images," in *CVPR*, 2016, pp. 2764–2772.
- [126] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [127] Y. Watkins, O. Iaroshenko, M. Sayeh, and G. Kenyon, "Image compression: Sparse coding vs. bottleneck autoencoders," in *IEEE Southwest Symposium on Image Analysis and Interpretation*. IEEE, 2018, pp. 17–20.
- [128] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the H.264/AVC video coding standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 560–576, 2003.
- [129] I. H. Witten, R. M. Neal, and J. G. Cleary, "Arithmetic coding for data compression," *Communications of the ACM*, vol. 30, no. 6, pp. 520–541, 1987.
- [130] D. H. Wolpert and W. G. Macready, "No free lunch theorems for optimization," *IEEE Transactions on Evolutionary Computation*, vol. 1, no. 1, pp. 67–82, 1997.
- [131] C.-Y. Wu, N. Singhal, and P. Krähenbühl, "Video compression through image interpolation," in *ECCV*, 2018, pp. 416–431.
- [132] F. Wu, D. Liu *et al.*, "Description of SDR video coding technology proposal by University of Science and Technology of China, Peking University, Harbin Institute of Technology, and Wuhan University," *JVET*, Tech. Rep. JVET-J0032, 2018.
- [133] J. Xu, M. Xu, Y. Wei, Z. Wang, and Z. Guan, "Fast H.264 to HEVC transcoding: A deep learning method," *IEEE Transactions on Multimedia*, DOI: 10.1109/TMM.2018.2885921, 2018.
- [134] M. Xu, T. Li, Z. Wang, X. Deng, R. Yang, and Z. Guan, "Reducing complexity of HEVC: A deep learning approach," *IEEE Transactions on Image Processing*, vol. 27, no. 10, pp. 5044–5059, 2018.
- [135] N. Yan, D. Liu, B. Li, H. Li, T. Xu, and F. Wu, "Convolutional neural network-based invertible half-pixel interpolation filter for video coding," in *ICIP*, 2018, pp. 201–205.
- [136] N. Yan, D. Liu, H. Li, B. Li, L. Li, and F. Wu, "Convolutional neural network-based fractional-pixel motion compensation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 3, pp. 840–853, 2019.
- [137] N. Yan, D. Liu, H. Li, and F. Wu, "A convolutional neural network approach for half-pel interpolation in video coding," in *ISCAS*. IEEE, 2017, pp. 1–4.
- [138] R. Yang, M. Xu, T. Liu, Z. Wang, and Z. Guan, "Enhancing quality for HEVC compressed videos," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2018.2867568, 2018.
- [139] R. Yang, M. Xu, Z. Wang, and T. Li, "Multi-frame quality enhancement for compressed video," in *CVPR*, 2018, pp. 6664–6673.
- [140] K. Yu, C. Dong, L. Lin, and C. C. Loy, "Crafting a toolchain for image restoration by deep reinforcement learning," in *CVPR*, 2018, pp. 2443–2452.
- [141] H. Zhang, L. Song, Z. Luo, and X. Yang, "Learning a convolutional neural network for fractional interpolation in HEVC inter coding," in *VCIP*. IEEE, 2017, pp. 1–4.
- [142] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [143] Q. Zhang, D. Liu, and H. Li, "Deep network-based image coding for simultaneous compression and retrieval," in *ICIP*. IEEE, 2017, pp. 405–409.
- [144] Y. Zhang, L. Sun, C. Yan, X. Ji, and Q. Dai, "Adaptive residual networks for high-quality image restoration," *IEEE Transactions on Image Processing*, vol. 27, no. 7, pp. 3150–3163, 2018.
- [145] Y. Zhang, T. Shen, X. Ji, Y. Zhang, R. Xiong, and Q. Dai, "Residual highway convolutional neural networks for in-loop filtering in HEVC," *IEEE Transactions on Image Processing*, vol. 27, no. 8, pp. 3827–3841, 2018.
- [146] Z. Zhang, Z. Chen, J. Lin, and W. Li, "Learned scalable image compression with bidirectional context disentanglement network," *arXiv preprint arXiv:1812.09443*, 2018.
- [147] L. Zhao, H. Bai, A. Wang, and Y. Zhao, "Learning a virtual codec based on deep convolutional neural network to compress image," *arXiv preprint arXiv:1712.05969*, 2017.
- [148] Z. Zhao, S. Wang, S. Wang, X. Zhang, S. Ma, and J. Yang, "Enhanced bi-prediction with convolutional neural network for high efficiency video coding," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2018.2876399, 2018.
- [149] L. Zhou, C. Cai, Y. Gao, S. Su, and J. Wu, "Variational autoencoder for low bit-rate image compression," in *CVPR Workshops*, 2018, pp. 2617–2620.