

# Neural Video Coding Using Multiscale Motion Compensation and Spatiotemporal Context Model

Haojie Liu<sup>ID</sup>, Ming Lu<sup>ID</sup>, Graduate Student Member, IEEE, Zhan Ma<sup>ID</sup>, Senior Member, IEEE, Fan Wang, Zhihuang Xie, Xun Cao<sup>ID</sup>, Member, IEEE, and Yao Wang<sup>ID</sup>, Fellow, IEEE

**Abstract**—Over the past two decades, traditional block-based video coding has made remarkable progress and spawned a series of well-known standards such as MPEG-4, H.264/AVC and H.265/HEVC. On the other hand, deep neural networks (DNNs) have shown their powerful capacity for visual content understanding, feature extraction and compact representation. Some previous works have explored the learnt video coding algorithms in an end-to-end manner, which show the great potential compared with traditional methods. In this paper, we propose an end-to-end deep neural video coding framework (NVC), which uses variational autoencoders (VAEs) with joint spatial and temporal prior aggregation (PA) to exploit the correlations in intra-frame pixels, inter-frame motions and inter-frame compensation residuals, respectively. Novel features of NVC include: 1) To estimate and compensate motion over a large range of magnitudes, we propose an unsupervised multiscale motion compensation network (MS-MCN) together with a pyramid decoder in the VAE for coding motion features that generates multiscale flow fields, 2) we design a novel adaptive spatiotemporal context model for efficient entropy coding for motion information, 3) we adopt nonlocal attention modules (NLAM) at the bottlenecks of the VAEs for implicit adaptive feature extraction and activation, leveraging its high transformation capacity and unequal weighting with joint global and local information, and 4) we introduce multi-module optimization and a multi-frame training strategy to minimize the temporal error propagation among P-frames. NVC is evaluated for the low-delay causal settings and compared with H.265/HEVC, H.264/AVC and the other learnt video compression methods following the common test conditions, demonstrating consistent gains across all popular test sequences for both PSNR and MS-SSIM distortion metrics.

**Index Terms**—Neural video coding, neural network, multiscale motion compensation, pyramid decoder, multiscale compressed flows, nonlocal attention, spatiotemporal priors, temporal error propagation.

Manuscript received July 9, 2020; revised September 22, 2020; accepted October 27, 2020. Date of publication November 3, 2020; date of current version August 4, 2021. This work was supported in part by the National Natural Science Foundation of China under Grant 62022038, in part by the China Scholarship Council under Grant 201906190086, and in part by OPPO Research Fund. This article was recommended by Associate Editor H. Meng. (*Corresponding authors:* Zhan Ma; Xun Cao.)

Haojie Liu, Ming Lu, Zhan Ma, and Xun Cao are with the School of Electronic Science and Engineering, Nanjing University, Nanjing 210093, China (e-mail: haojie@smail.nju.edu.cn; luming@smail.nju.edu.cn; mazhan@nju.edu.cn; caoxun@nju.edu.cn).

Fan Wang and Zhihuang Xie are with OPPO, Inc., Nanjing 210093, China (e-mail: wangfan6@oppo.com; xiezhihuang@oppo.com).

Yao Wang is with the Faculty of Electrical and Computer Engineering, Tandon School of Engineering, New York University, New York, NY 11201 USA (e-mail: yw523@nyu.edu).

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCSVT.2020.3035680

## I. INTRODUCTION

COMPRESSED video, a dominant media representation across the entire Internet, occupies more than 70% total traffic volume nowadays for entertainment (e.g., YouTube), productivity (e.g., tele-education), security (e.g., surveillance) etc. It still keeps growing explosively. Thus, in pursuit of efficient storage and network transmission, and pristine quality of experience (QoE) with higher resolution content (e.g., 2K, 4K and even 8K video with frame rate at 30 Hz, 60 Hz or even more), a better compression approach is greatly and continuously desired. In principle, the key problem in video coding is how to efficiently exploit visual signal redundancy using prior information, spatially (e.g., intra prediction, transform), temporally (e.g., inter prediction), and statistically (e.g., entropy context adaptation) for more compact representations with less bit rate consumption at the same reconstruction quality. This is well formulated as the minimization of Lagrangian cost  $J$  of rate-distortion optimization (RDO) that is widely adopted in existing video coders, e.g.,

$$\min J = R + \lambda \cdot D, \quad (1)$$

with  $R$  and  $D$  represent the compressed bit rates and reconstructed distortion respectively.

### A. Motivation

Over the past three decades, video compression technologies have been evolving and adapting constantly with coding efficiency improved by several folds, mostly driven under the efforts from the experts in ISO/IEC Moving Picture Experts Group (MPEG), ITU-T Video Coding Experts Group (VCEG) and their joint task forces. It leads to several popular video coding standards, including the H.264/Advanced Video Coding (H.264/AVC) [1], High-Efficiency Video Coding (HEVC) [2] and emerging versatile video coding (VVC) [3]. These standards share the similar (recursive) block-based hybrid prediction/transform framework where individual coding tools, such as the intra/inter prediction, integer transforms, context-adaptive entropy coding, etc, are intensively handcrafted to optimize the overall efficiency. Among them, *pixel-domain predictive coding* is one of the most important factors, contributing to the major performance gains [4]. For example, pixel-domain intra prediction was officially adopted into the H.264/AVC and later extended with the support of recursive block-sizes and abundant predictive directions for

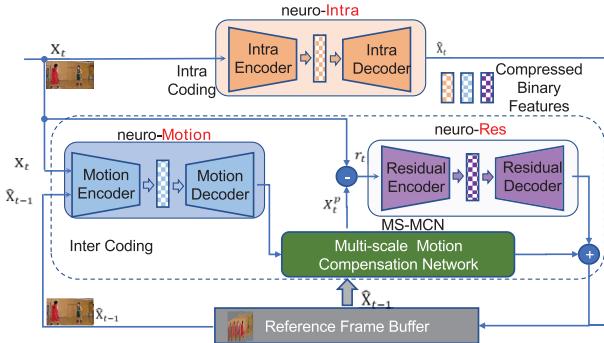


Fig. 1. Neural Video Coding (NVC). The modules neuro-Intra, neuro-Res, and neuro-Motion follow the general model architecture in Fig. 2 for efficient representations of intra pixels, displaced inter residuals, and inter motions. The neuro-Motion uses a pyramid decoder for the main decoder as discussed in Sec. III-D2.

efficiently exploiting spatial structures; recursive and even non-squared blocks are extensively used in inter prediction to remove temporal coherency. Basically, conventional video coding methods leverage the spatiotemporal pixel neighbors (as well as their linear combinations) for predictive signal construction, resulting in corresponding residuals for subsequent transform, quantization, and entropy coding for more compact representation. Optimal coding mode with appropriate block size and orientation (e.g., intra direction, inter motion vectors) is selected via computational RDO process, utilizing  $\ell_1$ -norm (e.g., mean absolute error - MAE) or  $\ell_2$ -norm (e.g., mean squared error - MSE) as the distortion metric.

Though recursive block-based pixel prediction shows its great success, it is mainly due to the hardware advancements in past decades, by which we can exhaustively search for the best prediction. It, however, is more and more challenging to simply trade computational resources for efficiency improvement because Moore's Law does not hold any more [5]. It therefore calls for innovative methodologies and architectures of video coding to further improve the coding efficiency, in response to the ever increasing users' requirement for video resolution and quality. Such pixel prediction strategy, either intra or inter, mostly relies on the physical coherence of video signal and applies the mathematical tools (e.g., linear weighting, orthogonal transform, Lagrangian optimization) for signal energy compaction.

### B. Our Approach

We propose an end-to-end neural video coding framework (NVC), which codes intra-frame pixels (called neuro-Intra), inter-frame motion (called neuro-Motion), and inter-frame residual (called neuro-Res) using separate variational autoencoders (VAEs), as shown in Fig. 1. A multiscale motion compensation network (MS-MCN) works together with neuro-Motion to generate multiscale optical flows and perform multiscale motion-compensated prediction of the current frame from the previous frame. The sparse image differences between past and present frame, e.g., residuals, are then encoded to obtain the final reconstruction; All three VAEs,

e.g., neuro-Intra, neuro-Motion, neuro-Res for compressing intra-pixel, inter-motion and inter-residual, are engineered together with MS-MCN in an end-to-end learning manner. Note that neuro-Intra takes a native image frame as input, neuro-Motion uses the current and past reconstructed frame, generating multiscale compressed flows (MCFs), MS-MCN uses these generated MCFs for motion compensation to obtain the inter-predicted frame and neuro-Res encodes the difference between the current frame and its prediction for the final reconstruction. Additionally, joint spatiotemporal and hyper priors are aggregated for efficient and adaptive context modeling of latent features to improve entropy coding efficiency for the motion field.

We have evaluated the efficiency of the proposed NVC for the low-delay causal settings against well-known HEVC, H.264/AVC and other learnt video compression methods following the common test conditions. The NVC demonstrated the leading performance with consistent gains across all popular test sequences for both PSNR (Peak signal-to-noise ratio) and MS-SSIM (multiscale structural similarity) [6] distortion metrics. Using the H.264/AVC as a common anchor, our NVC presents 35% BD-Rate (Bjontegaard Delta Rate) [7] gains, while HEVC and DVC (Deep Video Coding) [8] offer 30% and 22% gains, respectively, when the distortion is measured in terms of PSNR. Gains are even larger, if the distortion metric is replaced by the MS-SSIM. In this case, NVC can achieve 50% improvement, while both HEVC and DVC are around 25%. We further compare our NVC with DVC\_Pro [9] (an upgraded version of DVC) on HEVC test sequences: our NVC reveals 32.20% and 50.07% BD-Rate reduction measured by PSNR and MS-SSIM distortion respectively, while DVC\_Pro gives 34.57% and 45.88%.

Ablation studies have also been conducted to examine the gains due to different modules of NVC. We have shown that temporal priors and multi-frame training could greatly improve the efficiency and learning stability. Our MS-MCN is able to remove motion compensation noise by multiscale compensation, even better than using a cascaded trained denoising network.

### C. Novelty

The main contributions are highlighted below:

- We propose an end-to-end deep neural video coding framework (NVC), leveraging learnt feature domain representations for intra-pixel, inter-motion and inter-residual, respectively for compression;
- neuro-Motion and multiscale motion compensation network (MS-MCN) are employed together to capture coarse-to-fine motion displacements and obtain the prediction by warping the features of the reference frame at multiple scales;
- We propose a novel spatiotemporal context modeling approach for the entropy coding of the motion features, where the temporal context is obtained through modeling the temporal evolution of the motion features using a ConvLSTM in the temporal updating module (TUM). The temporal context is combined with the autoregressive

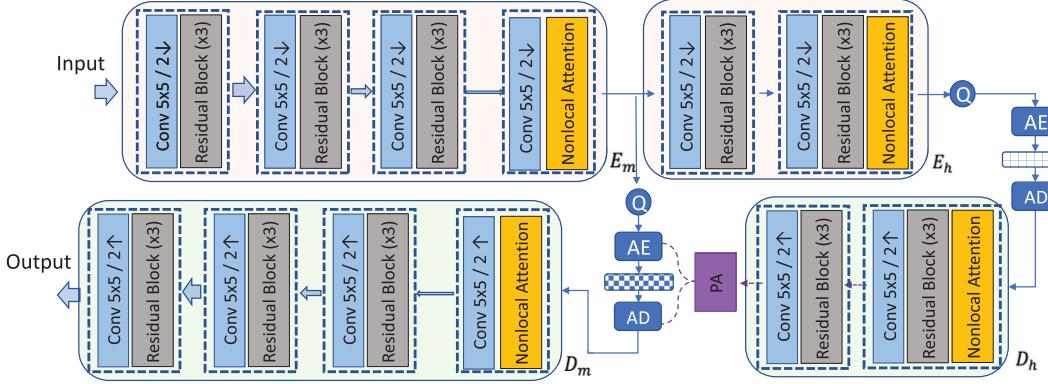


Fig. 2. The general structure of the variational autoencoder based compression engine used for neuro-Intra, neuro-Motion and neuro-Res in Fig. 1. The main encoder  $E_m$  includes major convolutional layers (each including a convolutional layer followed by down-sampling and three residual blocks, and the last layer includes a non-local attention module). The hyper encoder  $E_h$  includes two major convolutional layers. The main decoder  $D_m$  and hyper decoder  $D_h$  reverse the processing of  $E_m$  and  $E_h$ , respectively. Prior aggregation (PA) engine collects the information from hyper prior, autoregressive spatial neighbors, as well as temporal priors (if applicable) for efficient modeling of the probability distribution of latent features generated by the main encoder. Nonlocal attention is adopted at the bottlenecks of both main and hyper encoders to enable saliency based bit allocation, and rectified linear unit (ReLU) is embedded with convolutions for enabling the nonlinearity. “Q” is for quantization, AE and AD are for arithmetic encoding and decoding, respectively.  $2\downarrow$  and  $2\uparrow$  are downsampling and upsampling at a factor of 2 for both horizontal and vertical dimensions.

spatial context and hyperprior features in a spatiotemporal hyper aggregation module (STHAM);

- Nonlocal attention is attached at bottleneck layers of the VAE modules for adaptive bits allocation based on joint global and local feature extraction implicitly.

This work is based on our preliminary work [10] but with significant extensions and discussions including pyramid flow decoder in neuro-Motion, replacing single scale motion compensation with multiscale motion compensation, spatiotemporal prior aggregation engine for context modeling and progressive training with multiple frames.

The rest of this paper is structured as follows: Sec. II will briefly review relevant studies about learned image and video coding. Our neural video coding framework (NVC) is given in Sec. III with detailed discussions about neuro-Intra, neuro-Motion, MS-MCN and neuro-Res; Sec. IV will present experiments and ablation studies; and concluding remarks are drawn in Sec. V.

## II. RELATED WORK

Built on advancements of deep neural networks (DNNs), we have seen the explosive growth of DNN-based image/video compression approaches. Some explorations have attempted to replace modular components in traditional image/video coding framework such as filtering, intra prediction, etc; and some others have fully relied on powerful learning tools to perform the end-to-end optimization. Given that our neural video coding framework (NVC) belongs to the end-to-end learning category, we will emphasize the reviews in this avenue. For the modular optimization using DNNs, two great review articles can be found in [11], [12].

### A. Learnt Image Compression

DNN-based image compressions usually utilize auto-encoder or VAE architectures, consisting of nonlinear transforms, attention or importance map, differentiable

quantization, context model, and embedded loss functions, for end-to-end learning.

1) *Recurrent Autoencoder*: Toderic *et al.* [13] first proposed to utilize fully-connected recurrent auto-encoders for variable-rate thumbnail image compression. A serial improvements were then extended, including the full-resolution image support, learned entropy coding, unequal bits allocation, etc [14], [15] by the introductions of ConvLSTM or ConvGRU, for better coding efficiency. Variable bit rate is intrinsically enabled by such recurrent structure. It, however, suffers from the higher computational complexity at higher bit rates, because more recurrent processing are desired.

2) *Convolutional Autoencoder*: Alternatively, convolutional auto-encoders [16]–[19], are extensively studied in past years where different bit rates are adapted by setting a variety of  $\lambda$  in learning to optimize (1). Note that different network models may be required for individual bit rates, making it difficult for hardware implementation (e.g., model adaptation for various bit rates). Recently, conditional convolution [20] and scaling factor [21] were proposed to enable variable-rate compression using a single or a limited number of network parameters without noticeable coding efficiency loss. It makes the convolutional autoencoders more attractive for practical applications.

3) *Attention/Importance Map*: Li *et al.* [22] utilized a separate three-layer CNN to generate importance map for spatial-complexity-based adaptive bits allocation, leading to the noticeable subject quality improvement with well-preserved edges and textures. Instead, Mentzer *et al.* [18] further selected one channel from the bottleneck layer to unequally weigh features at different spatial locations for simplification. Such importance map embedding was lightweight and easy for training and end-to-end optimization. This approach was later improved with nonlocal attention mechanism to efficiently and implicitly capture both global and local important information for better compression [21].

4) *Nonlinear Transform*: To generate more compact feature representation, Balle *et al.* [16] suggested to replace the

traditional nonlinear activation, e.g., ReLU, with the generalized divisive normalization (GDN) that was theoretically proven to be more consistent with image natural statistics for visual perception. A succeeding investigation by the same authors was given in [23], reporting that GDN outperformed other nonlinear rectifiers, such as ReLU, leakyReLU and tanh, in compression tasks. Several follow-up studies [24], [25] directly applied GDN in their networks for compression exploration.

5) *Adaptive Contexts*: Probabilistic model plays a vital role in data compression. Assuming the Gaussian distribution for feature elements, Balle *et al.* [17] utilized hyper priors to estimate the parameters of Gaussian Scale Model (GSM) for latent features. Later Hu *et al.* [26] used hierarchical hyper priors (coarse-to-fine) for improving the entropy models in multiscale representations. Minnen *et al.* [19] improved the context modeling using joint autoregressive spatial neighbors and hyper priors based on Gaussian Mixture Model (GMM). Autoregressive spatial priors were usually extracted by PixelCNNs or PixelRNNs [27], which have been widely adopted for natural image density modeling. Reed *et al.* [28] further introduced multiscale PixelCNNs, yielding competitive density estimation and great speedup (e.g., from  $O(N)$  to  $O(\log N)$ ). It was later extended from 2D architectures to 3D PixelCNNs [18]. Channel-wise weights sharing-based 3D implementations can greatly reduce network parameters with higher parallelization. Then Chen *et al.* [21] discussed parallel pipelines of 3D PixelCNNs for practical decoding. Previous methods accumulated all the accessible priors to estimate the probability based on a single Gaussian distribution for each element. Recent explorations have shown that weighted GMMs can further improve the coding efficiency as reported by [29], [30].

6) *Quantization*: Quantization is a non-differentiable operation, basically converting continuous variables into discrete variables with a limited alphabet. This process has to be replaced by a differentiable operation when used in end-to-end learning framework for back propagation. A number of methods, such as uniform noise adding [16], stochastic rounding [13], soft-to-hard vector quantization [18] and universal quantization [20], were developed to approximate a continuous distribution for differentiation.

7) *Loss Functions*: Pixel-error, such as MSE, or MAE, was one of the most popular loss functions used. In the meantime, SSIM or MS-SSIM was also adopted because of its better consistency with visual perception. Simulations had revealed that SSIM-based loss can improve the perception quality, especially at low bit rates. Towards the perception-optimized encoding, perceptual losses that were measured by adversarial loss [31]–[33] and VGG loss [34] were embedded in learning to produce visually appealing results.

### B. Learnt Video Compression

Learnt video compression is extended from the learnt image compression by further exploiting the temporal redundancy in a trainable way for efficient representations of temporal motion and displaced image difference (e.g., predictive residual).

In most cases, network models for image compression were directly re-used for temporal displaced residuals, leaving a great deal of efforts devoted for better motion field compression.

Chen *et al.* [35] developed the DeepCoder where a simple convolutional autoencoder was applied for both intra and residual coding at fixed  $32 \times 32$  blocks, and block-based motion estimation in traditional video coding was re-used for temporal compensation. Lu *et al.* [8] introduced the optical flow for motion representation in their DVC work, which, together with the intra coding in [17], demonstrated similar performance compared with the HEVC. However, the coding efficiency suffer from a sharp loss at low bit rates. Liu *et al.* [10] extended their nonlocal attention optimized image compression (NLAIC) for coding the intra and residual frame, and applied spatiotemporal adaptive context model for more compact motion representation, showing consistent rate-distortion gains across different contents and bits rates.

Motion can be also implicitly inferred by temporal interpolations. For example, Wu *et al.* [36] applied recurrent neural network (RNN)-based frame interpolation. Together with the residual compensation, it offered comparable performance with H.264/AVC. Djelouah *et al.* [37] further imporved the interpolation-based video coding by utilizing the advanced optical flow estimation and feature domain residual coding. Yang *et al.* [38] also use a interpolation method with three hierarchical quality layers and a recurrent enhancement network for compression. However, temporal interpolation usually led to coding delay that may not be acceptable for low-latency applications.

Another interesting exploration made by Ripple *et al.* in [39] was to jointly encode the flow and residual signals using unified quantized features in an unsupervised way. A recurrent state was embedded to aggregate multi-frame information for efficient flow generation and residual coding.

## III. NEURAL VIDEO CODING

### A. Overview of NVC

Our NVC framework is designed for low-delay applications. As with all modern video encoders, the proposed NVC compresses the first frame in each group of pictures as an intra-frame using a VAE based compression engine (neuro-Intra). It codes remaining frames using motion compensated prediction. As shown in Fig. 1, it uses the VAE compressor (neuro-Motion) to generate the multiscale motion field between the current frame and the reference frame. Then, MS-MCN takes multiscale compressed flows, warps the multiscale features of the reference frame, and combines these warped features to generate the predicted frame. The prediction residual is then coded using another VAE-based compressor (neuro-Res).

Given a group of pictures (GOP)  $\mathbb{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t\}$ , we first encode  $\mathbf{X}_1$  using the neuro-Intra module, leading to the reconstructed frame  $\hat{\mathbf{X}}_1$ . The following frame  $\mathbf{X}_2$  is encoded predictively, using neuro-Motion, MS-MCN, and neuro-Res together, shown in Fig. 1. Note that MS-MCN takes the multiscale optical flows  $\{\vec{f}_d^1, \vec{f}_d^2, \dots, \vec{f}_d^s\}$  derived by the pyramid

TABLE I  
ABBREVIATIONS AND NOTATIONS

abbr.	description
NLAM	NonLocal Attention Module
LAM	Local Attention Module
MS-MCN	Multi-scale Motion Compensation Network
SS-MCN	Single-scale Motion Compensation Network
PA	Prior Aggregation
MCF	Multiscale Compressed Flow
neuro-Intra	Neural intra coding
neuro-Motion	Neural motion coding
neuro-Res	Neural residual coding
MSE	Mean Squared Error
MAE	Mean Absolute Error
PSNR	Peak signal-to-noise ratio
MS-SSIM	Multiscale Structural Similarity

decoder in neuro-Motion, and then generates the predicted frame  $\hat{\mathbf{X}}_2^p$  by multiscale motion compensation. Displaced inter-residual  $\mathbf{r}_2 = \mathbf{X}_2 - \hat{\mathbf{X}}_2^p$  is then compressed in neuro-Res, yielding the reconstruction  $\hat{\mathbf{r}}_2$ . The final reconstruction  $\hat{\mathbf{X}}_2$  is given by  $\hat{\mathbf{X}}_2 = \hat{\mathbf{X}}_2^p + \hat{\mathbf{r}}_2$ . Encoding of the next P-frame follows the same procedure of  $\mathbf{X}_2$  until all frames in the GOP are coded completely.

Table I summarizes relevant abbreviations used throughout this paper.

### B. The VAE Architecture Using NLAM and Spatiotemporal Priors

The general architecture of the VAE model is shown in Fig. 2, with main encoder-decoder pair for latent feature analysis and synthesis, and hyper encoder-decoder for hyper prior generation. The main encoder  $E_m$  uses four stacked CNN layers, where each convolutional layer employs stride convolutions to achieve downsampling (e.g., at a factor of 2 in this example) and cascaded convolutions (e.g., three ResNet-based residual blocks [40]<sup>1</sup>) for efficient feature extraction. We utilize two-layer hyper encoder  $E_h$  to further generate the subsequent hyper priors as side information, which are used for the entropy coding for the latent features.

To capture the spatial locality, we apply convolutional layers with limited receptive field (e.g.,  $3 \times 3$ ) that are stacked altogether to simulate the layer-wise feature extraction. These same ideas are used in many relevant studies [17], [19]. We utilize the simplest ReLU as the nonlinear activation function. Other nonlinear activation functions can be used as well, such as the GDN (Generalized Divisive Normalization) in [16].

Attention mechanism is leveraged to intelligently allocate bit resource (e.g., via unequal feature quantization) for image/video compression [18], [41]. It basically enables more accurate reconstruction of salient areas. We adopt the nonlocal attention module (e.g., NLAM) at the bottleneck layers of both the main encoder and hyper encoder, prior to quantization to include both global and local information for more accurate

<sup>1</sup>We apply cascaded ResNets for stacked CNNs because of its high-efficiency and reliable performance. Other efficient CNNs architectures can be applied as well.

importance selection. This module is motivated by the behaviour of HVS, where we often promptly scan the entire viewing scene to have the complete understanding of the field of vision, and then fixate to the salient regions.

To enable more accurate conditional probability density modeling for entropy coding of the latent features, we introduce the *Prior Aggregation* (PA) engine which fuses the inputs from the hyper priors, spatial neighbors and temporal context (if applicable).<sup>2</sup> The more accurate context modeling requires less resource (e.g., bits) for information representation as suggested in information theory [42]. For the sake of simplicity, we assume the latent features (e.g., motion, image pixel, residual) following the Gaussian distribution as in [19], [26], and use the PA engine to derive the mean and standard deviation of the distribution for each feature.

### C. Neural Intra Coding

Our neuro-Intra is a simplified version of the NLAIC that was originally proposed in [21].

1) *NLAIC*: One major distinction of NLAIC from the VAE model using autoregressive spatial context in [19] is the introduction of a nonlocal attention module (NLAM) inspired by [43]. NLAM is used to capture the global and local importance for saliency aggregation, by which we can assign different importance to various spatial-channel elements implicitly. Such nonlocal attention mechanism is inspired by the visual information processing of spatial perception (e.g., coarse global structure plus fine-grain local details). We have found that with the addition of NLAM, we can achieve similar performance as [19] without using GDN nonlinearity.

In addition, we have applied 3D  $5 \times 5 \times 5$  masked CNN<sup>3</sup> to extract spatial priors, which are fused with hyper priors in PA for entropy context modeling (e.g., bottom part of Fig. 6). Here, we have assumed the single Gaussian distribution for the context modeling of entropy coding. More details of NLAIC can be found in [21]. Note that temporal priors are not adopted for intra-pixel and inter-residual in this paper.

2) *Improvements to NLAIC*: Original NLAIC applies multiple NLAMs in both main and hyper coders, leading to excessive memory consumption at a large spatial scale. In NVC, NLAMs are only used at the bottleneck layers for both main and hyper encoder-decoder pairs, achieving adaptive bits allocation implicitly.

To overcome the non-differentiability of the quantization operation, quantization is simulated by adding uniform noise in [17]. However, such noise augmentation is not exactly consistent with the rounding in inference, yielding the performance loss as reported by [20]. Thus, we apply the universal quantization (UQ) [20] in neuro-Intra, i.e.:

$$\hat{x} = \mathbb{R}(x + u) - u, \quad (2)$$

where  $\hat{x}$  is a quantized symbol and  $u$  represents the random uniform variable ranging from  $-\frac{1}{2}$  to  $\frac{1}{2}$ . Statistically, we define

<sup>2</sup>Intra and residual coding only use joint spatial and hyper priors without temporal inference.

<sup>3</sup>This  $5 \times 5 \times 5$  convolutional kernel shares the same parameters for all channels, offering great model complexity reduction compared with 2D CNN-based solution in [19].

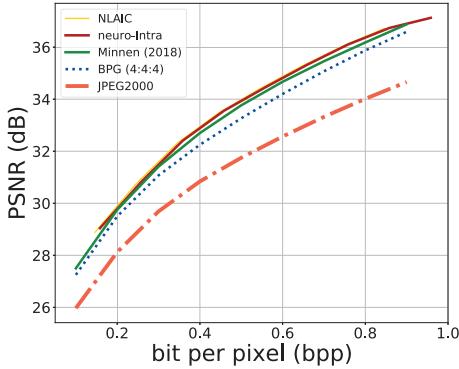


Fig. 3. **Efficiency of neuro-Intra.** PSNR vs. rate performance of neuro-Intra in comparison to NLAIC [21], Minnen (2018) [19], BPG (4:4:4) and JPEG2000. Note that the curves for neuro-Intra and NLAIC overlap.

the gradient as 1 for back propagation since UQ can be approximated as a linear function. Such UQ (2) is used for neuro-Motion and neuro-Res as well. When applying to common Kodak dataset, neuro-Intra achieved similar performance as NLAIC [21], outperforming Minnen (2018) [19], BPG (4:4:4) and JPEG2000, as shown in Fig. 3.

#### D. Neural Motion Coding and Compensation

Inter-frame coding plays a vital role in video coding. The key is how to efficiently represent motion in a compact format for effective compensation. In comparison to the pixel-domain block-based motion estimation and compensation in conventional video coding, we rely on the optical flow to accurately capture the temporal information for *motion compensation*.

Single scale pixel-level motion compensation can hardly handle the problem of occlusions and large motions. There are two main advantages from using multiscale motion estimation and compensation. Firstly, the multiscale motion estimation can capture motion fields from different scales which can better solve the problem of large motion. Meanwhile, the feature-level motion compensation at each scale can achieve feature-level inter prediction at each generated scale. Then the prediction will be achieved by a coarse-to-fine refinement architecture which can better handle the occlusion problems with progressive non-linear processing, compared with the single-scale linear warping in pixel domain.

1) *Single-Scale Motion Generation and Compensation:* Deep learning has spawned several optical flow estimation methods in either supervised or unsupervised way, such as FlowNet2 [44], PWC-net [45], etc. They mostly focus on acquiring high-precision optical flow between two consecutive (uncompressed) frames without any compression rate constraint. In video coder, it, however, is more challenging to derive robust flow where bitrate is often limited, and the reference frame is usually lossy encoded with inherent compression noises.

One approach applies a two-stage method shown in Fig. 4a, which is utilized by DVC [8]. Such two-stage scheme first uses a pre-trained lossless flow generator (e.g., FlowNet2) for explicit flow derivation, and then cascades an auto-encoder to encode the flow. RDO is examined when compressing the

optical flow to balance the bits allocation and reconstruction distortion. Either separable or joint optimization of flow derivation and compression can be applied for such two-stage approach.

We propose an alternative one-stage framework, shown in Fig. 4b, which is first presented in [10]. It directly transforms concatenated two frames (e.g., one is the reference from past, and one is current frame) into quantized temporal features that represent the inter-frame motion. These quantized features are decoded into compressed optical flow in an unsupervised way for frame compensation via warping. Such one-stage scheme does not require any pre-trained flow network such as FlowNet2 or PWC-net to generate the optical flow explicitly. It allows us to quantize the motion features rather than the optical flows and train the motion feature encoder and decoder together with explicit consideration of quantization and rate constraint. Note that the motion features are generated from the main motion encoder, and hyper encoder is used to generate hyper features for motion features in the VAE model in Fig. 2.

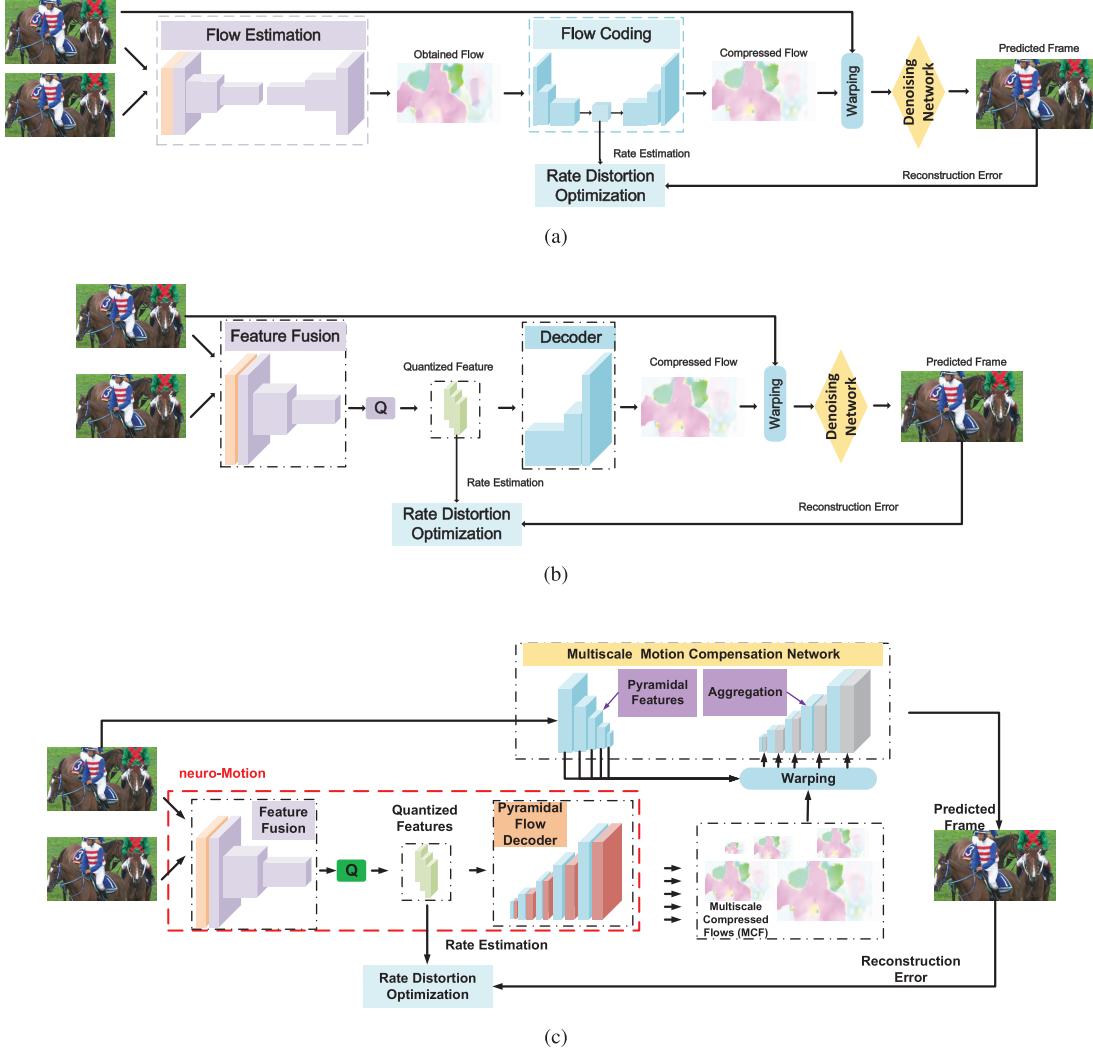
Loop filters, or quality enhancement networks, can be used in both two-stage and one-stage methods to enhance the quality of flow-warped frame, for improved quality of predicted frame, as shown in Fig. 4a and 4b. This is mainly because linear interpolation-based backward warping would inevitably introduce artifacts in reconstruction, especially for cases having flow estimation errors [46]. A cascaded quality enhancement model can be devised to alleviate such artifacts and improve the reconstruction quality. End-to-end learning can be further applied to the entire framework.

2) *Multiscale Motion Generation and Compensation:* We further extend our earlier work [10] to multiscale motion generation and compensation, for better inter-frame prediction. The neuro-Motion is modified for multiscale motion generation, where the main encoder is used for feature fusion, but we replace the main decoder by a *pyramidal flow decoder*, which generates the Multiscale Compressed optical Flows (MCFs). MCFs will be processed together with the reference frame, using a *multiscale motion compensation* network (MS-MCN) to obtain the predicted frame efficiently, as shown in Fig. 4c. The overall pipeline has four key steps:

- **Step 0: Pretrain the neuro-Motion using  $\ell_1$  loss for MCFs generation.** This part involves the implicit motion feature derivation in encoder, and multiscale (compressed) flow decoding in decoder.

For the former part, we use the main encoder of the VAE (see Fig. 2) to extract temporal features for implicit motion representation. It inputs the concatenation of reference frame  $\hat{\mathbf{X}}_{t-1}$  and current frame  $\mathbf{X}_t$  to derive latent motion features. To ensure the accurate MCFs generation, we replace the residual blocks in both VAE encoder and decoder with ResNet-based local attention modules (LAMs).<sup>4</sup> Here, quantized features have a size of  $(H/16) \times (W/16) \times 192$  with  $H, W$  denoting respective height and width of the original frame.

<sup>4</sup>Local attention module only utilizes local features to generate the attention maps compared with NLAM.



**Fig. 4. Motion estimation and compensation.** (a) Two-stage single-scale motion coding and compensation approach using a pre-trained flow network (with explicit raw flow) and a cascaded flow compression network (e.g., [8]); (b) One-stage unsupervised motion generation and compensation approach with implicit flow represented by quantized features that will be decoded into a single scale motion field for motion compensation (warping) (e.g., [10]); (c) One-stage neuro-Motion with MS-MCN uses a pyramidal flow decoder to synthesize the multiscale compressed flows (MCFs) that are used in a multiscale motion compensation network for generating predicted frames.

For the pyramidal flow decoder, we plug additional convolutional layers at each scale after the LAM to produce resolution-dependent flows  $\vec{f}_d^s, s = 0, 1, \dots, 4$  to capture the multiscale motion fields. The pyramid flow decoder is first trained in an unsupervised way using the loss:

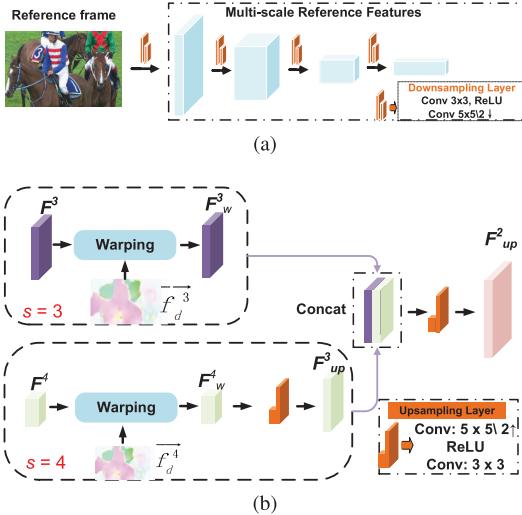
$$L = \sum_{s=0}^{n=4} \alpha_s \cdot \mathbb{D}_1(\hat{\mathbf{X}}_t^s, \mathbf{X}_t^s), \quad (3)$$

with  $\mathbb{D}_1$  as  $\ell_1$  loss.  $\alpha_s$ s are scale-dependent weighting coefficients which is set as  $\alpha_s = 4^s$  empirically. Here, we apply warping loss at each scale to train the optical flow generation and then warp the optical flow to the different scales of features generated from the reference frames to get the warped features at each scale.  $s = 0$  refers to the finest resolution.  $\hat{\mathbf{X}}_t^s$  is obtained by backward warping in each scale:

$$\hat{\mathbf{X}}_t^s = \text{warping}(\hat{\mathbf{X}}_{t-1}^s, \vec{f}_d^s). \quad (4)$$

Multiscale labels  $\mathbf{X}_t^s$ , and multiscale references  $\hat{\mathbf{X}}_{t-1}^s$  are generated from the current and reference frames respectively using average pooling with stride 2 for both vertical and horizontal dimensions.

- **Step 1: Fine-tune neuro-Motion via rate-constrained  $\ell_1$  loss.** We further fine-tune the neuro-Motion with rate constraints, where the bit rate of quantized feature representations will be estimated based on adaptive contexts generated from the hyper features and other spatiotemporal context (to be discussed in Fig. 6). We use a rate-distortion loss as in (1) where the distortion is the multiscale prediction loss as in (3) and the rate is the estimated entropy of the motion features.
- **Step 2: Pretrain MS-MCN via MCFs.** Using the MCFs (e.g.,  $\vec{f}_d^s$ s) generated by the neuro-Motion from **Step 1**, we will pretrain the *multiscale motion compensation network* for motion compensated prediction in the feature domain. It first uses a pyramidal decomposition of the



**Fig. 5. Multiscale motion compensation.** (a) A pyramidal decomposition is applied on reference frame to generate multi-scale features. It consecutively uses a  $3 \times 3$  convolution, ReLU activation and a  $5 \times 5$  convolution with stride 2. (b) Feature Aggregation between consecutive scales. Each upsampling layer consists of a  $5 \times 5$  convolution with stride 2, ReLU activation and a  $3 \times 3$  convolution.

reference frame, shown in Fig. 5a, to generate multiscale features  $\{\mathbf{F}^s\}$ ,  $s = 0, 1 \dots 4$  of the reference frames  $\hat{\mathbf{X}}_{t-1}$  which are respectively warped with the MCFs at corresponding scale for progressive aggregation.

Fig. 5b exemplifies the aggregation at the smallest with  $s = 4$  and the second smallest scale with  $s = 3$ : we use  $f_d^4$  to warp corresponding  $\mathbf{F}^4$  using (4), leading to the warped representation  $\mathbf{F}_w^4$ . This  $\mathbf{F}_w^4$  is then upsampled to  $\mathbf{F}_{up}^3 = \mathbb{U}(\mathbf{F}_w^4)$  that has the same resolution as the next scale, and concatenated with  $\mathbf{F}_w^3$ . These concatenated features are then upsampled again to yield scale 2 features:

$$\mathbf{F}_{up}^2 = \mathbb{U}\left(CAT(\mathbf{F}_w^3, \mathbf{F}_{up}^3)\right) = \mathbb{U}\left(CAT(\mathbf{F}_w^3, \mathbb{U}(\mathbf{F}_w^4))\right). \quad (5)$$

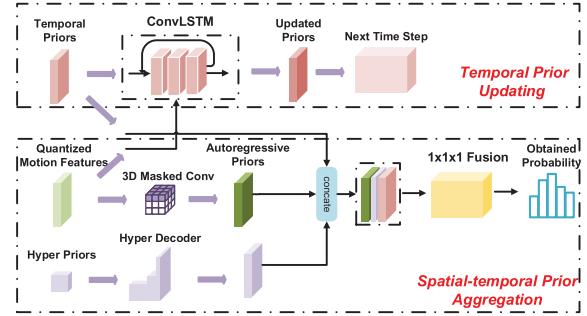
Note that  $\mathbf{F}_{up}^2$  will be concatenated with  $\mathbf{F}_w^2$  using similar steps. The upsampling operator  $\mathbb{U}()$  is implemented with two convolution layers as detailed in Fig. 5b. Eventually, we will apply a fusion layer, consisting of a  $1 \times 1$  convolution, ReLU and a  $3 \times 3$  convolution, on  $CAT(\mathbf{F}_w^0, \mathbf{F}_{up}^0)$  to obtain the final predicted frame. Pretraining the MS-MCN is accomplished by minimizing the multiscale motion compensation loss in (3).

- **Step 3: End-to-end Joint Refinement.** In the end, we take pre-trained neuro-Motion and MS-MCN together and perform an end-to-end joint refinement. We use the rate-constrained loss as in **Step 1**, and use the reconstructed intra frame as the reference frame. We resort to the “pre-training and joint-refinement” strategy since applying the joint training directly makes the model extremely unstable and hard to converge as revealed experimentally.

Table II reports the consistent PSNR gains offered by the MS-MCN in Sec. III-D2, in comparison to the SS-MCN [10]

TABLE II  
COMPARATIVE STUDIES OF MS-MCN AND SS-MCN USING THE PSNRs OF PREDICTED FRAME

Sequences	High Bit Rate		Low Bit Rate	
	MS-MCN	SS-MCN	MS-MCN	SS-MCN
BasketballPass	<b>32.94</b>	29.34	<b>30.52</b>	30.01
RaceHorses	<b>27.48</b>	25.08	<b>26.12</b>	24.41
PartyScene	<b>27.01</b>	26.43	<b>25.55</b>	25.27
BQMall	<b>33.09</b>	31.35	<b>30.76</b>	30.33
vidyo1	<b>38.81</b>	37.14	<b>36.34</b>	36.30
vidyo4	<b>38.86</b>	36.60	<b>36.40</b>	36.33
Average	<b>33.03</b>	30.99	<b>30.94</b>	30.44



**Fig. 6. Context-adaptive modeling using joint spatio-temporal and hyper priors.** All priors are fused in PA to provide estimates of the probability distribution parameters.

in Sec. III-D1, revealing the efficiency of MS-MCN for accurately motion compensated prediction. At high bit rate scenario, almost 2 dB gain is observed; MS-MCN keeps superior efficiency even with challenging temporal motions or occlusions in “BasketballPass” content. At low bit rates, gains are reduced but still noticeable. It shows that the PNSR gains depend on how much motion information are actually compressed. The more information (e.g., high bit rate) comes with larger PSNR gains, and vice versa.

3) *Context-Adaptive Flow Coding*: It is well recognized that motion fields present high correlations, both spatially and temporally. To exploit such correlation, Ripple *et al.* [39] proposed to code the flow differences in feature domain. Here, we propose to exploit this correlation in the context modeling for the entropy coding of the motion features. Specifically, we develop a joint spatiotemporal and hyper prior-based context-adaptive model shown in Fig. 6. This is implemented in PA engine of Fig. 2 for neuro-Motion.

The proposed PA engine for neuro-Motion consists of a *spatio-temporal-hyper aggregation module* (STHAM) and a *temporal updating module* (TUM), shown in Fig. 6. At timestamp  $t$ , STHAM accumulates all the accessible priors and estimate the mean and standard deviation of the assumed Gaussian distribution for each new quantized motion feature  $\hat{\mathcal{F}}$  using:

$$(\mu_{\hat{\mathcal{F}}}, \sigma_{\hat{\mathcal{F}}}) = \mathbb{G}(\hat{\mathcal{F}}_1, \dots, \hat{\mathcal{F}}_{t-1}, \hat{\mathbf{z}}_t, \mathbf{h}_{t-1}), \quad (6)$$

Here,  $\hat{\mathcal{F}}_i, i = 1, 2, \dots$  are elements of quantized latent motion features for the current frame,  $\mathbf{h}_{t-1}$  is consists of temporal priors derived from the motion features preceding the current frame.  $\hat{\mathbf{z}}_t$  includes hyper priors of the quantized

motion features. These features are concatenated and fused using stacked  $1 \times 1 \times 1$  convolutions. Note that masked convolution is used for the spatial features  $\hat{\mathcal{F}}$ .

To generate the temporal priors, TUM is applied to current quantized features  $\hat{\mathcal{F}}_t$  recurrently using a standard ConvLSTM:

$$(\mathbf{h}_t, \mathbf{c}_t) = \text{ConvLSTM}(\hat{\mathcal{F}}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}), \quad (7)$$

where  $\mathbf{h}_t$  are updated temporal priors for the next frame,  $\mathbf{c}_t$  is a memory state to control information flow across multiple time instances (e.g., frames). Note that the STHAM and TUM components are trained with other components in neuro-Motion in **Step 1** and **Step 3** described in Sec. III-D2.

It is worth to point out that leveraging temporal correlation for compact motion representation is also widely explored in traditional video coding approaches. For example, motion vector predictions from spatial and temporal co-located neighbors are standardized in HEVC, by which only motion vector differences (after prediction) are encoded. Here, instead of coding the flow feature differences, we use the flow features in the past to help estimating the probability distribution of the flow features in the current frame.

#### E. Neural Residual Coding

Inter-frame residual coding is another significant module contributing to the overall system efficiency. It is used to compress the temporal prediction error pixels. It affects the efficiency for next frame prediction since errors usually propagate temporally.

Here we use the VAE architecture in Fig. 2 for encoding the residual  $\mathbf{r}_t$ . The rate-constrained loss is used:

$$L = \lambda \cdot \mathbb{D}_2(\mathbf{X}_t, (\mathbf{X}_t^P + \hat{\mathbf{r}}_t)) + R, \quad (8)$$

where  $\mathbb{D}_2$  is the  $\ell_2$  loss between a residual compensated frame  $\mathbf{X}_t^P + \hat{\mathbf{r}}_t$  and  $\mathbf{X}_t$ . neuro-Res will be first pretrained using the predicted frames by the pretrained neuro-Motion and MS-MCN, and a loss function in (8) where the rate  $R$  only consider the bits for residual. Then we refine it jointly with neuro-Motion and MS-MCN, using a loss where  $R$  considers the bits for both motion and residual with two frames.

#### F. Training Strategy

**1) Progressive Training:** Training all NVC components directly is difficult and unreliable, because these modules are interdependent with each other. For example, inter prediction depends on the former reconstructed frame and meanwhile a better predicted frame will reduce the residual energy. In our model development, we first train the neuro-Intra for multiple bit rates, using multiple  $\lambda$  values; Then we pretrain and fine-tune the neuro-Motion and MS-MCN through **Step 1-3** described in Sec. III-D2, to minimize a training loss that consider both the prediction error and the rate for motion features. We also pretrain neuro-Res using the predicted frames generated by the trained neuro-Motion and MS-MCN so far, with a loss that considers both bit rate of residual and the final reconstruction of the current frame. Finally, we further refine all the modules for inter-coding, including

neuro-Motion, MS-MCN and neuro-Res, using a training loss including the reconstruction error and the total rate for the motion features and residual features. For each target bit rate, we set the  $\lambda$  for training the inter-coding model to be proportional to the  $\lambda$  used for training the neuro-Intra module.

**2) Training With Multi-Frame Loss:** One way to train the inter-coding model (including neuro-Motion, MS-MCN and neuro-Res modules) is to use only two frames as a training sample and use the neuro-Intra module to code the first frame, and use the inter-coding model to code the second frame. This training approach can only learn the intra-to-inter variations, yielding instability and poor reconstruction for the succeeding inter frames distant from the first intra frame. To overcome the quality degradation in future inter-frames, we adopt a multi-frame training strategy.

We first pretrain the inter-coding model using pairs of two successive frames as training samples, where the first frame is encoded and decoded using the neuro-Intra. This step follows the progressive training procedure. We then refine the inter-coding model by using groups of four successive frames as training samples, where the first frame is encoded and decoded by a pretrained neuro-Intra, and each following frame is coded using the inter-coding model and the decoded frame is then used recursively as the reference frame for coding the next frame. We use a loss function that considers the sum of the distortions (in MSE or negative MS-SSIM) in all three P-frames, and the sum of the rate for these frames, so that the model update considers the impact of the propagation of P-frame reconstruction errors. Note that once the inter-coding model is trained, it is applied to all inter frames in a GOP during testing. Although it is possible to improve the performance by training a different inter-coding model for each possible distance between a P-frame and the intra-frame, we choose to train a single model that is used repeatedly for all P-frames, to reduce the implementation complexity.

## IV. EXPERIMENTAL STUDIES

#### A. Datasets and Hyperparameters for Model Training

**1) Datasets:** The neuro-Intra is trained on COCO [47] and CLIC [48] with training samples randomly cropped into  $256 \times 256 \times 3$ . The neuro-Motion, MS-MCN and neuro-Res are joint trained using Vimeo 90k [49] with sample size of  $192 \times 192 \times 3$ .

Our NVC is evaluated using the standard HEVC test sequences and ultra video group (UVG) dataset. HEVC test sequences include different classes, covering the contents with a variety of motion, frame rate, resolution and texture; UVG dataset has seven 1080p videos which are often selected for testing video applications.

**2) Loss Function:** We have used either MSE or negative MS-SSIM loss for training the intra-coding and inter-coding modules. For pretraining the neuro-Motion and MS-MCN modules, we use the  $\ell_1$  loss on the prediction error instead which is similar to the mean absolute error (MAE) used in traditional motion estimation. We do not use MS-SSIM loss on the predicted frame because MS-SSIM mainly cares about

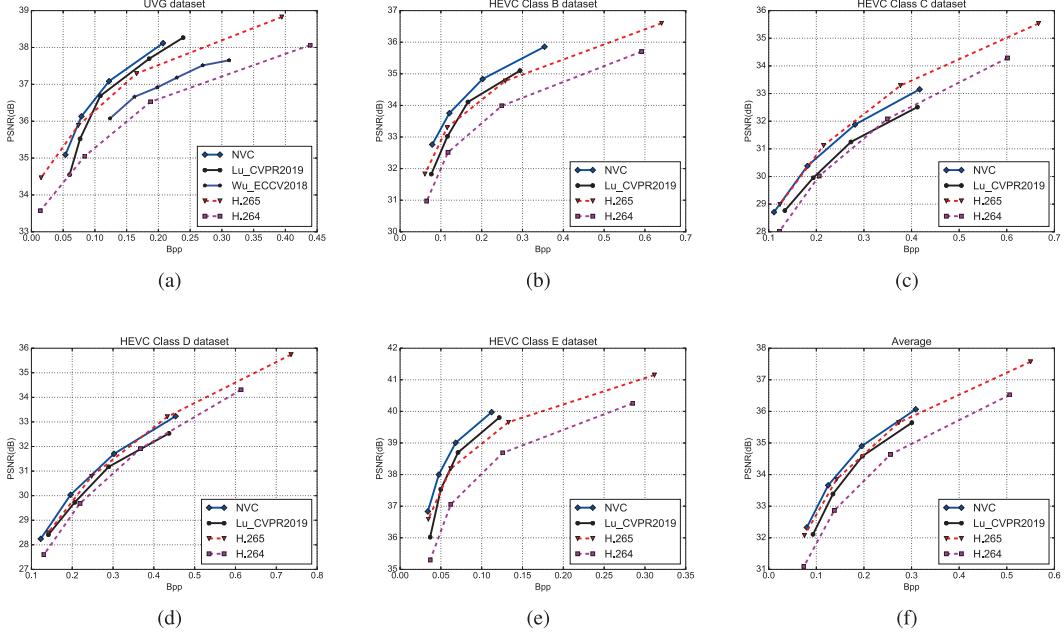


Fig. 7. PSNR vs. rate performance.

TABLE III  
BD-RATE GAINS OF NVC, HEVC AND DVC AGAINST THE H.264/AVC

Sequences	H.265/HEVC				DVC				NVC			
	PSNR		MS-SSIM		PSNR		MS-SSIM		PSNR		MS-SSIM	
	BDBR	BD-(D)	BDBR	BD-(D)	BDBR	BD-(D)	BDBR	BD-(D)	BDBR	BD-(D)	BDBR	BD-(D)
ClassB	-32.03%	0.78	-27.67%	0.0046	-27.92%	0.72	-22.56%	0.0049	<b>-45.66%</b>	<b>1.21</b>	<b>-54.90%</b>	<b>0.0114</b>
ClassC	<b>-20.88%</b>	<b>0.91</b>	-19.57%	0.0054	-3.53%	0.13	-24.89%	0.0081	-17.82%	0.73	<b>-43.11%</b>	<b>0.0133</b>
ClassD	-12.39%	0.57	-9.68%	0.0023	-6.20%	0.26	-22.44%	0.0067	<b>-15.53%</b>	<b>0.70</b>	<b>-43.64%</b>	<b>0.0123</b>
ClassE	-36.45%	0.99	-30.82%	0.0018	-35.94%	1.17	-29.08%	0.0027	<b>-49.81%</b>	<b>1.70</b>	<b>-58.63%</b>	<b>0.0048</b>
UVG	-48.53%	1.00	-37.5%	0.0056	-37.74%	1.00	-16.46%	0.0032	<b>-48.91%</b>	<b>1.24</b>	<b>-53.87%</b>	<b>0.0100</b>
Average	-30.05%	0.85	-25.04%	0.0039	-22.26%	0.65	-23.08%	0.0051	<b>-35.54%</b>	<b>1.11</b>	<b>-50.83%</b>	<b>0.0103</b>

the structure similarity and may ignore the background noise to some extent. Through experiments, we have found that using MS-SSIM loss for neuro-Motion can lead to temporal error accumulation, which not only increases the bits consumption but also potentially leads to unreliable training.

3) *Hyperparameters and Platform*: The initial learning rate (LR) is set to 10e-4 and is halved for every 10 epochs, and final models are obtained using a LR of 10e-5 for both pretraining and overall training. We apply the distributed training on 4 GPUs (Titan Xp) for 3 days for each bit rate model using deep learning framework PyTorch.

4) *Evaluation Criteria*: We apply the same low-delay coding setting as DVC in [8] for our method and traditional H.264/AVC, and HEVC for fair comparison. We encode 100 frames and use GOP of 10 on HEVC test sequences, and 600 frames with GOP of 12 on UVG dataset. Both PSNR and MS-SSIM results are offered to understand the efficiency of NVC.

5) *Generate Models for Different Target Rates*: To generate a model for a particular bit rate, we first train the neuro-Intra using a certain  $\lambda$  value, denoted by  $\lambda_{\text{intra}}$ . Then we train the inter-frame models (neuro-motion, MS-MCN and neuro-Res) using a proportionally reduced  $\lambda$  value,  $\lambda_{\text{inter}} = \lambda_{\text{intra}}/4$ . This

simple way of setting the  $\lambda_{\text{intra}}$  and  $\lambda_{\text{inter}}$  values yielded good results. We chose not to further optimize  $\lambda_{\text{inter}}$  for given  $\lambda_{\text{intra}}$ .

### B. Performance Comparison

1) *Rate Distortion Performance*: We compare the coding performance of different methods in Fig. 7 and 8 using respective PSNR and MS-SSIM measures, across HEVC and UVG test sequences. Note that when we report PSNR or MS-SSIM values, the models are trained using PSNR and MS-SSIM, respectively, as the distortion metric. In Table III, by setting the same anchor using H.264/AVC, our NVC presents 35% BD-Rate gains when the distortion is measured by PSNR, while HEVC offers 30% gains. Gains are even larger, if the distortion is measured by the MS-SSIM. In this case, NVC can achieve 50% improvement, while HEVC is around 25%.

Wu *et al.* [36] proposed a interpolation based video coding framework and could not get better performance than H.265/HEVC. Our NVC performs significantly better than DVC [8], which has 22.26% and 23.08% BD-Rate reductions under PSNR and MS-SSIM metrics, respectively. From Fig. 7 and Fig. 8, DVC [8] has mainly improved the coding efficiency against HEVC at high bit rates. However, DVC is not competitive at low bit rate (e.g., having performance

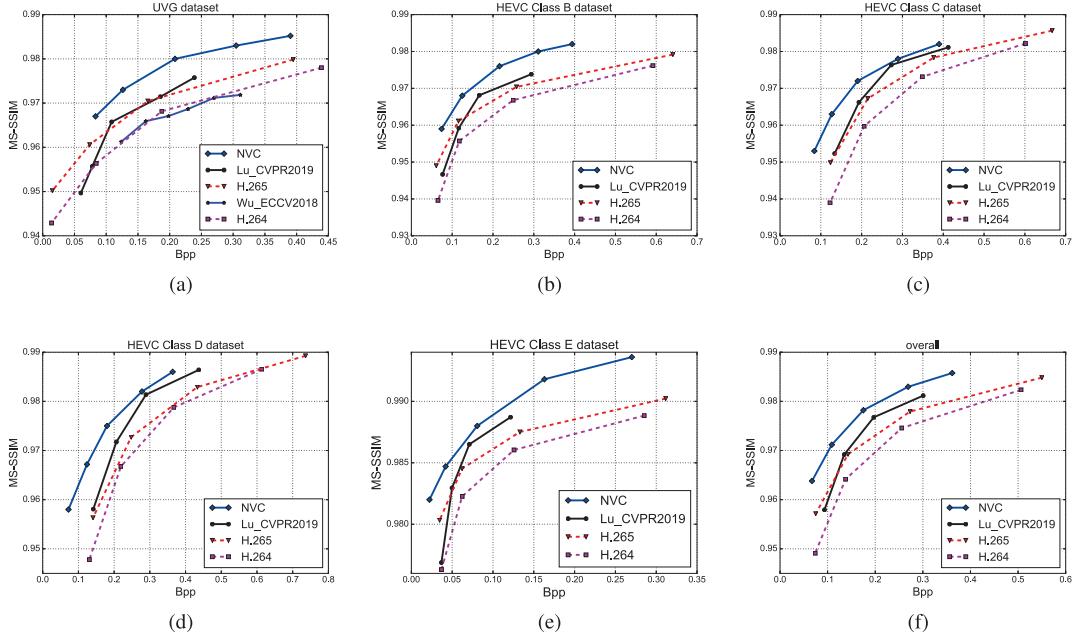


Fig. 8. **MS-SSIM vs. rate performance.** NVC shows significant gains for all the testing videos. MS-SSIM is usually more correlated with the perceptual quality than PSNR, especially at low bit rates.

even worse than H.264/AVC at some rate regimes.). We have also observed that DVC’s performance varies for different test sequences. An improved version of DVC, known as DVC\_pro, has been recently reported, which has shown the state of the art performance [9] by using [19] for intra and residual coding and  $\lambda$  tuning. NVC is slightly better than DVC Pro in both PSNR and MS-SSIM, at 34.5% and 45.88%, respectively.

2) *Visual Comparison*: We provide the visual quality comparison with H.264/AVC and H.265/HEVC in Fig. 10. Generally, NVC leads to a higher quality reconstruction at slightly lower bit rate. For example, for the “RaceHorse” which has nontranslational motion and complex background, NVC uses 7% less bits for more than 1.5 dB PSNR improvement compared with H.264/AVC. For other cases, our method also shows robust improvement. Traditional codec usually suffers from blocky artifacts and motion-induced noise close to object edges. One can clearly observe block partition boundaries with severe pixel discontinuity in reconstructed frames by H.264/AVC. Our results have much less visible noise and artifacts.

3) *Complexity Analysis*: We conduct several experiments on the Intel Xeon E5-2680 v4 CPU@2.40GHz and a TiTan XP GPU. We select two common official softwares JM and HM, two commercial softwares x264 and x265 and our proposed NVC for comparison. For a 1080p video sequence, the encoding time of JM and HM are 0.134fps and 0.0146fps, respectively. Our NVC provides 0.6fps encoding speed which is 41.09 times faster than HM. For x264 and x265, they are mainly optimized for practical and commercial use. Several complexity reduction and model simplification techniques are applied for the acceleration of the coding speed. The faster mode for the encoding can be nearly 112fps and 20fps for x264 and x265, respectively. With neural network acceleration

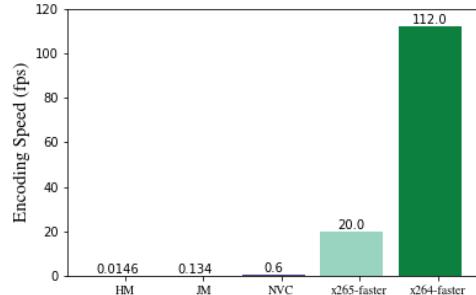


Fig. 9. **Encoding speed.**

and simplification such as network pruning, parallel processing, model quantization and further advances, NVC can be further simplified and improved for higher inference speed. Experiment results are included in Fig. 9.

### C. Ablation Study

This section examines modular components in NVC to further understand its capacity for application.

1) *Spatiotemporal Context Modeling*: To evaluate the gain from using temporal priors generated by ConvLSTM for the entropy coding of the motion features, we have also trained a model when the PA engine in neuro-Motion does not use the temporal priors. Table IV shows that 2% to 7% rate increase are incurred if the temporal priors are not used. This reveals that temporal priors help to make probability prediction more accurate, leading to less bit consumption for compressing the motion features.

Generally, bits consumed by motion information vary across different content and total target rates. More saving is reported for low bit rates, and for motion intensive content. In these



Fig. 10. **Visual comparison.** Reconstructed frames of NVC, HEVC and H.264/AVC. NVC has fewer blocky artifacts and visible noise, etc, and provides better quality at lower bit rate.

TABLE IV  
EFFICIENCY OF TEMPORAL PRIORS FOR CODING MOTION FEATURES. THE ENTRIES IN 2ND AND 3RD COLUMNS ARE THE BD-RATE REDUCTION

Sequences	NVC	W/O temporal priors	Bits saving
ClassB	-54.90%	-48.77%	6.13%
ClassC	-43.11%	-39.17%	3.94%
ClassD	-43.64%	-39.23%	4.41%
ClassE	-58.63%	-56.27%	2.36%
UVG dataset	-53.87%	-49.14%	4.73%
Average	-50.83%	-46.52%	4.31%

cases, motion bits usually occupy a higher percentage of the total rate. For stationary content, such as HEVC Class E, motion bits contribute less percentage, thus the gain from using temporal priors is less significant.

2) *Comparison With Cascaded Denoising Networks:* Linear warping often brings artifacts (e.g., ghosting edges), especially when flow estimation is not accurate and occlusion happens.

To remove such artifacts, a denoising network is often added after motion compensated warping as shown in Fig. 4a and 4b. Our MS-MCN, instead, rely on multiscale motion fields to generate the predicted frame. As a comparison, we also trained a U-net like denoising network on the warped frames based on the single-scale decoded motion field. As shown in Fig. 11, MS-MCN achieves 0.2 to 0.3 dB PSNR gains across a large bit rate ranges over using a cascaded denoising network for various content for overall performance.

3) *Temporal Stability With Multi-Frame Training:* Compression Errors (e.g., lossy compression noise) often propagate from one frame to another due to predictive coding structure. It is critical to limit the temporal quality degradation during training. To evaluate the gain from multi-frame training, we compare the performance of the NVC model trained by minimizing reconstruction errors for only one P-frame (single frame training) and the model further refined by minimizing the reconstruction loss for multiple consecutive

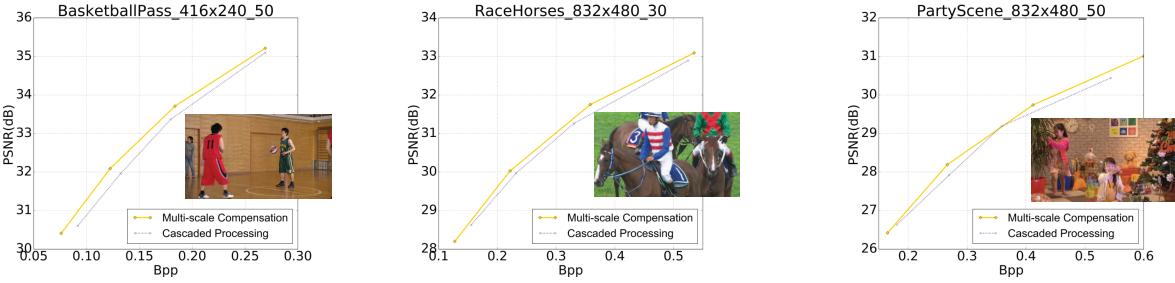


Fig. 11. **Comparison with using a cascaded denoising network.** Proposed MS-MCN can achieve better rate-distortion performance compared with using a denoising network following single scale motion compensation.

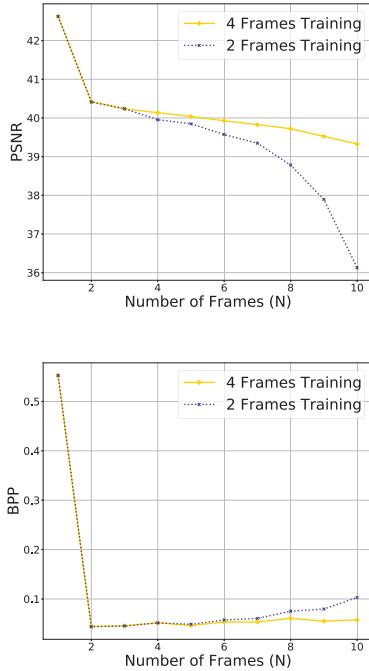


Fig. 12. **Benefit from multi-frame training.** Two-frame training could only learn temporal variations from intra to inter reconstruction; Multi-frame training is applied to capture both intra-to-inter, and inter-to-inter variations, yielding improved stability and performance.

P-frames. We have varied the number of total frames to use while training, and found that using 4 frames reaches a good compromise. As shown in Fig. 12, multi-frame training could well capture temporal quality variations from intra to inter, and from inter to inter, leading to a more generalized model with consistent performance for all the frames in a GOP with slower quality degradation, and improved stability.

## V. CONCLUSION

We have developed an end-to-end deep neural video coding framework (NVC) that can compactly represent the intra-pixel, inter-motion and inter-residual information, respectively. We have shown that the pyramid decoder in neuro-Motion and the multiscale motion compensation network (MS-MCN) together can significantly improve inter-frame prediction, compared to the more conventional single scale motion compensation. Furthermore, adding temporal context can lead to more

efficient entropy coding of the motion information than using only spatial context and hyper priors. We further demonstrate that using a multi-frame training loss can effectively mitigate the temporal error propagation.

We evaluate the NVC by PSNR and MS-SSIM respectively, and compare its performance both with standard coding methods including H.264/AVC and H.265/HEVC as well as learnt video coders including Wu *et al.* [36], DVC [8] and DVC\_pro [9]. NVC offers consistent and stable gains over existing methods across a variety of contents and bit rates.

The proposed DNN-based model can be further improved by engineering better stacked CNNs instead of current implementation. Context-adaptive probability models could be further improved. For example, recent exploration in [30] has demonstrated that weighted GMM could further improve the entropy coding efficiency. This can be easily incorporated in our PA engine. Reference frame selection is another direction for overall efficiency optimization, by which we can embed and aggregate most appropriate information for improving the inter-coding efficiency [50].

The H.264/AVC, HEVC, and even VVC, are engineering masterpieces for video coding.  $\lambda$  adaptation, rate control, etc can be definitely borrowed to improve the NVC. Furthermore, how to make NVC practically applicable is also worth for deep investigation.

To benefit the research community, all relevant materials will be made publically available soon at <https://njuvision.github.io/Neural-Video-Coding>.

## REFERENCES

- [1] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, “Overview of the H.264/AVC video coding standard,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 13, no. 7, pp. 560–576, Jul. 2003.
- [2] G. J. Sullivan, J.-R. Ohm, W.-J. Han, and T. Wiegand, “Overview of the high efficiency video coding (HEVC) standard,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 22, no. 12, pp. 1649–1668, Dec. 2012.
- [3] B. Bross, J. Chen, S. Liu, and Y.-K. Wang, *Versatile Video Coding (Draft 7)*, document JVET-P2001, Oct. 2019.
- [4] J.-R. Ohm, G. J. Sullivan, H. Schwarz, T. K. Tan, and T. Wiegand, “Comparison of the coding efficiency of video coding standards—Including high efficiency video coding (HEVC),” *IEEE Trans. circuits Syst. Video Technol.*, vol. 22, no. 12, pp. 1669–1684, Dec. 2012.
- [5] J. Dean, “1.1 the deep learning revolution and its implications for computer architecture and chip design,” in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2020, pp. 8–14.
- [6] Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multiscale structural similarity for image quality assessment,” in *Proc. 37th Asilomar Conf. Signals, Syst. Comput.*, vol. 2, 2003, pp. 1398–1402.

- [7] G. Bjontegaard, *Calculation of Average PSNR Differences Between RD-Curves*, document VCEG-M33, 2001.
- [8] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, “DVC: An end-to-end deep video compression framework,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 11006–11015.
- [9] G. Lu, X. Zhang, W. Ouyang, L. Chen, Z. Gao, and D. Xu, “An End-to-End learning framework for video compression,” *IEEE Trans. Pattern Anal. Mach. Intell.*, early access, Apr. 20, 2020, doi: 10.1109/TPAMI.2020.2988453.
- [10] H. Liu *et al.*, “Learned video compression via joint spatial-temporal correlation exploration,” in *Proc. AAAI*, 2020, pp. 1–8.
- [11] D. Liu, Y. Li, J. Lin, H. Li, and F. Wu, “Deep learning-based video coding: A review and a case study,” *ACM Comput. Surveys*, vol. 53, no. 1, pp. 1–35, May 2020.
- [12] S. Ma, X. Zhang, C. Jia, Z. Zhao, S. Wang, and S. Wang, “Image and video compression with neural networks: A review,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 6, pp. 1683–1698, Jun. 2020.
- [13] G. Toderici *et al.*, “Variable rate image compression with recurrent neural networks,” in *Proc. Int. Conf. Learn. Represent.*, 2016, pp. 1–12.
- [14] G. Toderici *et al.*, “Full resolution image compression with recurrent neural networks,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 5306–5314.
- [15] N. Johnston *et al.*, “Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 4385–4393.
- [16] J. Ballé, V. Laparra, and E. P. Simoncelli, “End-to-end optimized image compression,” 2016, *arXiv:1611.01704*. [Online]. Available: <http://arxiv.org/abs/1611.01704>
- [17] J. Ballé, D. Minnen, S. Singh, S. Jin Hwang, and N. Johnston, “Variational image compression with a scale hyperprior,” 2018, *arXiv:1802.01436*. [Online]. Available: <http://arxiv.org/abs/1802.01436>
- [18] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. V. Gool, “Conditional probability models for deep image compression,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, vol. 1, no. 2, pp. 4394–4402.
- [19] D. Minnen, J. Ballé, and G. D. Toderici, “Joint autoregressive and hierarchical priors for learned image compression,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 10794–10803.
- [20] Y. Choi, M. El-Khamy, and J. Lee, “Variable rate deep image compression with a conditional autoencoder,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 3146–3154.
- [21] T. Chen, H. Liu, Z. Ma, Q. Shen, X. Cao, and Y. Wang, “Neural image compression via non-local attention optimization and improved context modeling,” 2019, *arXiv:1910.06244*. [Online]. Available: <http://arxiv.org/abs/1910.06244>
- [22] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, “Learning convolutional networks for content-weighted image compression,” 2017, *arXiv:1703.10553*. [Online]. Available: <http://arxiv.org/abs/1703.10553>
- [23] J. Ballé, “Efficient nonlinear transforms for lossy image compression,” 2018, *arXiv:1802.00847*. [Online]. Available: <http://arxiv.org/abs/1802.00847>
- [24] J. Lee, S. Cho, and S.-K. Beack, “Context-adaptive entropy model for End-to-end optimized image compression,” 2018, *arXiv:1809.10452*. [Online]. Available: <http://arxiv.org/abs/1809.10452>
- [25] J. Klopp, Y.-C. F. Wang, S.-Y. Chien, and L.-G. Chen, “Learning a code-space predictor by exploiting intra-image-dependencies,” in *Proc. BMVC*, 2018, p. 124.
- [26] Y. Hu, W. Yang, and J. Liu, “Coarse-to-fine hyper-prior modeling for learned image compression,” in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 1–8.
- [27] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu, “Pixel recurrent neural networks,” 2016, *arXiv:1601.06759*. [Online]. Available: <http://arxiv.org/abs/1601.06759>
- [28] S. Reed *et al.*, “Parallel multiscale autoregressive density estimation,” in *Proc. 34th Int. Conf. Mach. Learning (JMLR)*, vol. 70, 2017, pp. 2912–2921.
- [29] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, “Learned image compression with discretized Gaussian mixture likelihoods and attention modules,” 2020, *arXiv:2001.01568*. [Online]. Available: <http://arxiv.org/abs/2001.01568>
- [30] J. Lee, S. Cho, and M. Kim, “An end-to-end joint learning scheme of image compression and quality enhancement with improved entropy minimization,” 2019, *arXiv:1912.12817*. [Online]. Available: <http://arxiv.org/abs/1912.12817>
- [31] O. Rippel and L. Bourdev, “Real-time adaptive image compression,” 2017, *arXiv:1705.05823*. [Online]. Available: <http://arxiv.org/abs/1705.05823>
- [32] C. Huang, H. Liu, T. Chen, Q. Shen, and Z. Ma, “Extreme image coding via multiscale autoencoders with generative adversarial optimization,” in *Proc. IEEE Vis. Commun. Image Process. (VCIP)*, Dec. 2019, pp. 1–4.
- [33] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. Van Gool, “Generative adversarial networks for extreme learned image compression,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 221–231.
- [34] H. Liu, T. Chen, Q. Shen, T. Yue, and Z. Ma, “Deep image compression via end-to-end learning,” in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, 2018, pp. 2575–2578.
- [35] T. Chen, H. Liu, Q. Shen, T. Yue, X. Cao, and Z. Ma, “DeepCoder: A deep neural network based video compression,” in *Proc. IEEE Vis. Commun. Image Process. (VCIP)*, Dec. 2017, pp. 1–4.
- [36] C.-Y. Wu, N. Singhal, and P. Krähenbühl, “Video compression through image interpolation,” in *Proc. ECCV*, 2018, pp. 416–431.
- [37] A. Djelouah, J. Campos, S. Schaub-Meyer, and C. Schroers, “Neural inter-frame compression for video coding,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 6421–6429.
- [38] R. Yang, F. Mentzer, L. Van Gool, and R. Timofte, “Learning for video compression with hierarchical quality and recurrent enhancement,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 6628–6637.
- [39] O. Rippel, S. Nair, C. Lew, S. Branson, A. Anderson, and L. Bourdev, “Learned video compression,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 3454–3463.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 770–778.
- [41] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, “Learning convolutional networks for content-weighted image compression,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 3214–3223.
- [42] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Hoboken, NJ, USA: Wiley, 2012.
- [43] Y. Zhang, K. Li, K. Li, B. Zhong, and Y. Fu, “Residual non-local attention networks for image restoration,” 2019, *arXiv:1903.10082*. [Online]. Available: <http://arxiv.org/abs/1903.10082>
- [44] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, “FlowNet 2.0: Evolution of optical flow estimation with deep networks,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 2462–2470.
- [45] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, “PWC-net: CNNs for optical flow using pyramid, warping, and cost volume,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 8934–8943.
- [46] M. Jaderberg *et al.*, “Spatial transformer networks,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 2017–2025.
- [47] T.-Y. Lin *et al.*, “Microsoft coco: Common objects in context,” in *Proc. Eur. Conf. Comput. Vis.* Cham, Switzerland: Springer, 2014, pp. 740–755.
- [48] (2018). *Challenge on Learned Image Compression*. [Online]. Available: <http://www.compression.cc/challenge/>
- [49] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman, “Video enhancement with task-oriented flow,” *Int. J. Comput. Vis.*, vol. 127, no. 8, pp. 1106–1125, Aug. 2019.
- [50] H. Li, B. Li, and J. Xu, “Rate-distortion optimized reference picture management for high efficiency video coding,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 22, no. 12, pp. 1844–1857, Dec. 2012.



**Haojie Liu** received the B.S. degree from Nanjing University, Nanjing, China, in 2016, where he is currently pursuing the Ph.D. degree with the School of Electronic Science and Engineering. His research interests include video communication and processing, machine learning, and computer vision.



**Ming Lu** (Graduate Student Member, IEEE) received the B.E. degree in electronic science and engineering from Nanjing University, China, in 2016, where he is currently pursuing the Ph.D. degree. His current research interests include image/video processing, including deep learning-based super resolution and image/video coding. He was a co-recipient of the 2018 ACM SIGCOMM Student Research Competition Finalist, and the 2020 IEEE MMSP Image Compression Grand Challenge Best Performing Solution.



**Zhan Ma** (Senior Member, IEEE) received the B.S. and M.S. degrees from the Huazhong University of Science and Technology (HUST), Wuhan, China, in 2004 and 2006, respectively, and the Ph.D. degree from New York University, New York, in 2011. From 2011 to 2014, he has been with Samsung Research America, Dallas, TX, USA, and Futurewei Technologies, Inc., Santa Clara, CA, USA. He is currently on the faculty of Electronic Science and Engineering School, Nanjing University, China. His current research interests include learning-based

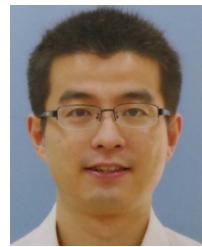
video coding and smart cameras. He was a co-recipient of the 2018 ACM SIGCOMM Student Research Competition Finalist, the 2018 PCM Best Paper Finalist, the 2019 IEEE Broadcast Technology Society Best Paper Award, and the 2020 IEEE MMSP Image Compression Grand Challenge Best Performing Solution.



**Fan Wang** received the B.S. degree from Nanjing University, Nanjing, China, in 2012, and the M.S. degree from Waseda University, Tokyo, Japan, in 2014. He was a Researcher with Samsung Research, Beijing, China, from 2017 to 2019. He is currently a Senior Researcher with OPPO Research, Beijing. His research interests include video coding and machine learning.



**Zhihuang Xie** received the M.S. degree in electronics engineering from Jinan University, Guangzhou, China, in 2018. From April 2018 to July 2018, he was a Research Assistant with the Department of Computer Science, City University of Hong Kong. He is currently working in video coding standardization as a Researcher with the Nebula Lab, OPPO Research Institute, Shenzhen, China. His research interests include video compression, machine learning, and super resolution.



**Xun Cao** (Member, IEEE) received the B.S. degree from Nanjing University, Nanjing, China, in 2006, and the Ph.D. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2012. He held visiting positions with Philips Research Aachen, Germany, in 2008, and Microsoft Research Asia, Beijing, from 2009 to 2010. He was a Visiting Scholar with The University of Texas at Austin, Austin, TX, USA, from 2010 to 2011. He is currently a Professor with the School of Electronic Science and Engineering, Nanjing University. His research interests include computational photography, image-based modeling and rendering, and VR/AR systems.



**Yao Wang** (Fellow, IEEE) received the B.S. and M.S. degrees in electronic engineering from Tsinghua University, Beijing, China, in 1983 and 1985, respectively, and the Ph.D. degree in electrical and computer engineering from the University of California at Santa Barbara in 1990.

Since 1990, she has been on the Faculty of Electrical and Computer Engineering, Tandon School of Engineering, New York University (formerly Polytechnic University, Brooklyn, NY). Her research areas include video communications, multimedia medical imaging. She is the leading author of a textbook titled *Video Processing and Communications*, and has published over 250 papers in journals and conference proceedings. She received the New York City Mayor's Award for Excellence in Science and Technology in the Young Investigator Category in 2000. She was an Elected Fellow of the IEEE in 2004, for contributions to video processing and communications. She is the Co-Winner of the IEEE Communications Society Leonard G. Abraham Prize Paper Award in the field of communications systems in 2004 and the Co-Winner of the IEEE Communications Society Multimedia Communication Technical Committee Best Paper Award in 2011. She has served as an Associate Editor for the IEEE TRANSACTIONS ON MULTIMEDIA and the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY.