# Frequency-aware Learned Image Compression for Quality Scalability

Hyomin Choi[†], Fabien Racapé[†], Shahab Hamidi-Rad[†], Mateen Ulhaq[*,†], Simon Feltman[†]

[†]Interdigital Emerging Technologies Lab, Los Altos, CA, USA

[*]School of Engineering Science, Simon Fraser University, Burnaby, BC, Canada

Email: {hyomin.choi; fabien.racape; shahab.hamidi-rad; mateen.ulhaq; simon.feltman}@interdigital.com

*Abstract*—Spatial frequency analysis and transforms serve a central role in most engineered image and video lossy codecs, but are rarely employed in neural network (NN)-based approaches. We propose a novel NN-based image coding framework that utilizes forward wavelet transforms to decompose the input signal by spatial frequency. Our encoder generates separate bitstreams for each latent representation of low and high frequencies. This enables our decoder to selectively decode bitstreams in a quality-scalable manner. Hence, the decoder can produce an enhanced image by using an enhancement bitstream in addition to the base bitstream. Furthermore, our method is able to enhance only a specific region of interest (ROI) by using a corresponding part of the enhancement latent representation. Our experiments demonstrate that the proposed method shows competitive rate-distortion performance compared to several non-scalable image codecs. We also showcase the effectiveness of our two-level quality scalability, as well as its practicality in ROI quality enhancement.

*Index Terms*—End-to-end compression, learned image compression, quality scalability, wavelet decomposition

## I. Introduction

Conventional lossy codecs [1] use transforms such as the Discrete Wavelet Transform (DWT) or Discrete Cosine Transform (DCT), alongside quantization to achieve variable-rate compression. An image is transformed into the frequency domain, and the resulting transformed coefficients are quantized into bins, where each bin is sized to minimize perceptible distortion. Most importantly, distortion at high spatial frequencies is much less noticeable than distortion at low frequencies. Based on this property of the human visual system, a quality-scalable coding method allows progressive improvement in decoded image quality by providing the decoder with further high-frequency information at higher bitrates [2]–[5].

End-to-end learned image compression (LIC) methods have recently caught the research community's interest. Ballé *et al.* [6] first proposed a density modeling method using Generalized Divisive Normalization (GDN) to *transform* the input images into an entropy coding-friendly latent space, which was effectively used in the autoencoder-based [7] approach in [8]. More recently, several variational autoencoder (VAE)-based methods [9]–[11] focused on accurately modeling the distributions of the latent variables, resulting in rate-distortion (RD) performance competitive with the latest fully-engineered codecs [12], [13]. Other approaches [14], [15] sought im-

provements in the analysis transform[1]. In particular, Akbari *et al.* [15] replaced regular 2D-convolutions with octave convolutions (OctConv) [16] which act like wavelet transforms in that the spatial resolution is reduced while diminishing spatial redundancy. However, rather than analyzing the input by frequency, the authors' modifications to OctConv focus on improving representational power.

We propose a quality-scalable frequency-aware learned image coding (FLIC) method using wavelet-embedded octave convolution (WeOctConv) that has RD performance competitive with non-scalable methods. Our method supports two-level quality scalability by encoding an input image into two separate bitstreams, as well as ROI-based quality scalability by encoding only selected regions of the latent space.

In Section II, we review prior work that inspired the proposed method, which is then described in detail in Section III. Experimental results are presented in Section IV, followed by conclusions in Section V.

## II. Prior work

Chen *et al.* [16] introduced the OctConv layer which factorizes its input into low-frequency ($L$) and high-frequency ($H$) features. Given an input tensor $\mathcal{Y}_{\text{in}} = \{\mathcal{Y}_{\text{in}}^H, \mathcal{Y}_{\text{in}}^L\}$, the output tensor $\mathcal{Y}_{\text{out}} = \{\mathcal{Y}_{\text{out}}^H, \mathcal{Y}_{\text{out}}^L\}$ is computed by

$$\begin{aligned}
\mathcal{Y}_{\text{out}}^H &= f^{H \to H}(\mathcal{Y}_{\text{in}}^H) + \texttt{Upsample}(f^{L \to H}(\mathcal{Y}_{\text{in}}^L)) \\
\mathcal{Y}_{\text{out}}^L &= f^{L \to L}(\mathcal{Y}_{\text{in}}^L) + f^{H \to L}(\texttt{Pool}(\mathcal{Y}_{\text{in}}^H))
\end{aligned} \quad (1)$$

where $f^{A \to B}$ represents the convolutional update between groups of frequencies $A$ and $B$. $f^{A \to B}$ is known as *inter-frequency communication* whenever $A$ and $B$ are distinct frequency groups, and *intra-frequency update* whenever they are the same. `Upsample` uses nearest-neighbor interpolation. `Pool` uses average pooling, which from the perspective of frequency analysis is similar to a low-pass filter for the 2-D discrete Haar wavelet transform. Thus, $\mathcal{Y}_{\text{out}}^L$ accumulates the low-frequency information over successive OctConv layers.

Akbari *et al.* [15] introduced a modified OctConv in which `Pool` is replaced with a convolution with a stride of 2. This modification increases representational ability, but is also less interpretable from the frequency decomposition perspective. Additionally, extensively applying GDN to the OctConv layers potentially adds redundant computations since OctConv already conducts some degree of factorization. Nonetheless, their approach shows improvements in RD performance in comparison to several LIC approaches.

---

[1]In LIC literature, the analysis transform $g_a$ transforms an input image into the latent space, from which the synthesis transform $g_s$ reconstructs the image.

Fig. 1. Design of the proposed layers (a) WeOctConv and (b) TWeOctConv, along with their corresponding Residual Block (RB). Inter-frequency updates are shown as purple lines, and intra-frequency updates as green lines.

## III. PROPOSED METHOD

This section describes the proposed WeOctConv layer and introduces our quality-scalable FLIC method, consisting of WeOctConvs, based on the factorized prior architecture [8].

### A. Wavelet-embedded Octave Convolution

Fig. 1 depicts the overall computation flow of the proposed WeOctConv layer, used in the analysis transform, and its dual TWeOctConv, used in the synthesis transform. The WeOctConv layer shown in Fig. 1(a) captures low and high frequency information from input tensors $\mathcal{Y}_{in}^{k} \in \mathbb{R}^{C_{in} \times N \times M}$ for each $k \in \{L, H\}$, where $C_{in}$ is the number of input channels and $N \times M$ is the feature resolution, and outputs the tensors $\mathcal{Y}_{out}^{k} \in \mathbb{R}^{C_{out} \times \frac{N}{2} \times \frac{M}{2}}$, where $C_{out}$ is the number of output channels. Conversely, the TWeOctConv layer in Fig. 1(b) synthesizes $\mathcal{Y}_{in}^{k} \in \mathbb{R}^{C_{in} \times \frac{N}{2} \times \frac{M}{2}}$ into $\mathcal{Y}_{out}^{k} \in \mathbb{R}^{C_{out} \times N \times M}$.

For the WeOctConv layer, the inter-frequency update consists of a DWT and a convolution, where the DWTs use the following $2 \times 2$ Haar wavelet kernels with a stride of 2:

$$LL = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad HH = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \quad (2)$$

The low-pass filter $LL$ is used for the $H \rightarrow L$ update, and the high-pass filter $HH$ is used for the $L \rightarrow H$ update. The downsampled filtered output is then fed into a convolutional layer with $3 \times 3$ kernels, denoted as Conv3. We found that it is challenging to simultaneously optimize for both the trainable DWT transform coefficients [17] and the RD criterion. Hence, we instead use fixed DWT transform coefficients. For the TWeOctConv layer, the natural choice for the inter-frequency update is a Conv3 and an inverse DWT (IDWT). However, we found that using fixed IDWT transform coefficients causes severe quality degradation, especially at lower bitrates. Hence, we instead use a Conv3PS layer composed of a Conv3 followed by a PixelShuffle [18], which is capable of upsampling its input in a way similar to IDWT by mixing groups of 4 channels.

Each intra-frequency update uses Residual Blocks (RBs). For WeOctConv, each RB consists of a Conv3 with stride 2 (denoted as Conv3s2), a L(eaky)ReLU, a Conv3, and a LReLU; and in the "skip" branch, a convolution with a $1 \times 1$ kernel and stride 2 (denoted as Conv1s2). For TWeOctConv, Conv3PS replaces both Conv3s2 and Conv1s2 within each RB.

Lastly, the inter-frequency updates are added to their corresponding intra-frequency updates.



Fig. 2. Overall architecture of the proposed FLIC. The bottom of each layer details a number of input channels $C_{in}$ and output channels $C_{out}$, respectively. For the first layer $C_{in}$ in the analysis and the last layer $C_{out}$ in the synthesis, only one port of input and output with three channels (i.e., RGB) is available.



Fig. 3. Tiled latent channels of (a) $\widehat{\mathcal{Y}}^{L}$ and (b) $\widehat{\mathcal{Y}}^{H}$ computed using a sample image from the Kodak dataset [20]. Each channel is normalised to $[0, 1]$ for visibility. All channels with non-zero latent representations are presented.

### B. Frequency-aware Learned Image Compression

Fig. 2 presents the overall architecture of the proposed image compression network employing (T)WeOctConvs. Denoted by "Frequency-aware Analysis", the encoder-side analysis transform $g_a(\cdot)$ with learned parameters $\psi_a$ analyzes an input image $\mathbf{X}$ and generates a pair of compressed latent representations:

$$\mathcal{Y} = \{\mathcal{Y}^{L}, \mathcal{Y}^{H}\} = g_a(\mathbf{X}; \psi_a). \quad (3)$$

We limit the use of GDNs to only one instance after the last WeOctConv. In our experiments, reducing the number of computationally intensive GDNs in this way does not result in a performance drop. Furthermore, no bias is used at all to reduce the number of operations. $\mathcal{Y}^{H}$ and $\mathcal{Y}^{L}$ are then rounded to the nearest integer during inference to obtain the quantized latent representations $\widehat{\mathcal{Y}}^{L}$ and $\widehat{\mathcal{Y}}^{H}$. During training, like in [8], uniform noise is added to the latent representations for the gradient computation. $\widehat{\mathcal{Y}}^{L}$ and $\widehat{\mathcal{Y}}^{H}$ are respectively losslessly encoded into a *base* and an *enhancement* bitstream using entropy encoders (ECs)[2].

The decoder uses entropy decoders (EDs) to losslessly decode the given bitstreams into $\widehat{\mathcal{Y}}^{L}$ and $\widehat{\mathcal{Y}}^{H}$. Denoted by "Frequency-aware Synthesis", the synthesis transform $g_s(\cdot)$ with learned parameters $\psi_s$ synthesizes a reconstructed image $\widehat{X}$ from $\widehat{\mathcal{Y}}^{L}$ and $\widehat{\mathcal{Y}}^{H}$ by applying to each a Conv3 and an inverse GDN (IGDN), and then applying a number of TWeOctConvs. This is written as

$$\widehat{\mathbf{X}} = g_s(\widehat{\mathcal{Y}}^{L}, \widehat{\mathcal{Y}}^{H}; \psi_s). \quad (4)$$

As shown in Fig. 3, spatial low- and high-frequency features are captured well by $\widehat{\mathcal{Y}}^{L}$ and $\widehat{\mathcal{Y}}^{H}$, respectively. This is because WeOctConv examines the wavelet decomposition on each layer's input throughout the analysis transform.

---

[2]Specifically, arithmetic range coder based on Asymmetric Numeral Systems (ANS) provided in [19] is used.

## C. Quality scalability

In traditional scalable codecs [2]–[4], the decoder may reconstruct an input image using its encoded base bitstream. A more detailed, higher quality image may be reconstructed by additionally providing the decoder with an enhancement bitstream that typically carries further high-frequency information. The proposed FLIC provides a similar form of quality scalability. Eq. (4) generates a high-quality input reconstruction using all latent representations. However, a lower quality image may be reconstructed by setting $\widehat{\mathcal{Y}}^H = \mathbf{0}$ so that

$$\widehat{\mathbf{X}}_{\text{base}} = g_s(\widehat{\mathcal{Y}}^L, \mathbf{0}; \psi_s), \tag{5}$$

where $\mathbf{0}$ denotes a zero tensor whose elements are all zeros with the same dimension as $\widehat{\mathcal{Y}}^H$.

Furthermore, the decoder also supports the quality enhancement of selected ROIs. This can be done by feeding $g_s$ with a decoded tensor $\widehat{\mathcal{Y}}^H_{\text{ROI}}$ (of the same dimensions as $\widehat{\mathcal{Y}}^H$) containing coded latent variables only for corresponding regions and zeros everywhere else. This produces the reconstructed image

$$\widehat{\mathbf{X}}_{\text{ROI}} = g_s(\widehat{\mathcal{Y}}^L, \widehat{\mathcal{Y}}^H_{\text{ROI}}; \psi_s), \tag{6}$$

where the selected ROI regions have been enhanced.

## D. Loss function

During training, we use a loss function in the form of an RD Lagrangian as in [8]:

$$\mathcal{L} = \underbrace{\mathbb{E}_{x \sim p_x}\left[-\log_2 p_{\hat{y}}(\hat{y})\right]}_{\text{rate estimation}} + \lambda \cdot \underbrace{(D(\mathbf{X}, \widehat{\mathbf{X}}) + \alpha \cdot D(\mathbf{X}, \widehat{\mathbf{X}}_{\text{base}}))}_{\text{distortion}}, \tag{7}$$

where $p_x$ denotes the probability density of the input data $x$ and $p_{\hat{y}}$ represents a fully factorized distribution of the quantized latent variable $\hat{y}$. The distortion metric $D$ can be any objective quality metric; we use mean squared error (MSE) and MS-SSIM [21] for our experiments. The hyperparameter $\alpha \geq 0$ balances the importance in quality of the full base+enhancement reconstruction $\widehat{\mathbf{X}}$ and the base-only reconstruction $\widehat{\mathbf{X}}_{\text{base}}$.

## IV. EXPERIMENTS

Our FLIC networks are trained on random cropped patches of size $256 \times 256$ from the Vimeo-90K dataset [22]. The mini-batch size is set to 8 and our networks are trained for up to 2.5M steps ($\approx 350$ epochs), corresponding to about 10 to 12 days. We use an Adam optimizer with an initial learning rate of $10^{-4}$, which is then decreased by 90% after the first 30 epochs whenever the validation loss plateaus with a patience of 4 epochs. We train models for each $\lambda = 2^n \cdot 10^{-2}$ over all $n \in \{3, 2, 1, 0, -1, -2, -3\}$. The models are tested on all 24 images in the Kodak dataset [20] to evaluate RD performance.

## A. Compression performance

*1) Impact of the hyperparameter $\alpha$:* We introduced $\alpha$ in Eq. (5) to control the quality of $\widehat{\mathbf{X}}_{\text{base}}$. To examine its impact on RD performance, we trained our model over various values of $\lambda$ for several $\alpha \in \{0.1, 0.01, 0.001, 0.0001, 0\}$. The RD curves for various $\alpha$ are shown in Fig. 4. The RD performance of the full bitstream (solid curves) is only marginally affected by



Fig. 4. RD performance of our proposed FLIC with various $\alpha$ on the Kodak dataset [20]. Solid and dashed lines represent the full bitstream and the base-only bitstream, respectively.



Fig. 5. Comparison of average bit proportion between the base layer and the enhancement layer for several models including [15] and the proposed models, trained with $\lambda = 2^3 \cdot 10^{-2}$.

changing $\alpha$. In contrast, the RD performance of the base-only bitstream (dashed curves) drastically improves as $\alpha$ increases. Thus, a larger value of $\alpha \geq 0.1$ may be used to obtain good base-only reconstructions without compromising full reconstruction performance.

Fig. 5 compares the average bit proportion between the base and enhancement bitstreams for various models. Although there is no scalability in [15][3], we consider their low- and high-frequency feature tensors to be the base and the enhancement layers, respectively. On average, about 4% of the entire bitstream by [15] represents low-frequency information. In contrast, our method adaptively uses $\alpha$ to control the bit proportion trade-off between bitstreams, without significantly affecting overall RD performance.

*2) Comparison with benchmarks:* Table I summarizes the average bit savings by several LIC methods and our proposed models with various $\alpha$ against JPEG2000[4] in terms of BD-rate [25]. Since our method is built upon [8] and uses 30M parameters, we reconfigured and retrained the relevant benchmarks [8], [15] to use the same number of parameters and entropy bottleneck as ours. As such, we reasonably compare the coding results of the benchmarks in the first two columns with our proposed methods in the last three columns of Table I. The first row shows the average bit savings of the models optimized for MSE in terms of BD-rate computed on the PSNR versus bpp curves. The last row shows coding gain by the models optimized for MS-SSIM in terms of BD-rate computed on the MS-SSIM versus bpp curves.

---

[3]Since there is no publicly available code, we have reimplemented the network using the factorized prior-based entropy bottleneck [8] and trained it on the same dataset used for our proposed model.

[4]FFMPEG [23] (v3.4.8) and libopenjpeg [24] (v2.3.0) are used.

Fig. 6. Visual example of our quality scalability method for the sample input image `kodim14.png` from the Kodak [20] dataset. The first column (a) represents the uncompressed input. The top header shows the coding results (RGB-PSNR / bpp) for the last three columns in the case of: (b) using both base and enhancement bitstreams, (c) using only the base bitstream, and (d) using the base bitstream along with the enhancement bitstream containing only selected ROIs. For visibility, the first row shows a cropped patch of the compressed image at (300, 50) with a size of $468 \times 400$. The second row shows the image represented in the Fourier domain. The third row shows enlarged ROIs.

TABLE I.    BD-RATE (%) PERFORMANCE OF VARIOUS LEARNED IMAGE CODECS AGAINST JPEG2000 [3]

| Anchor | Benchmark LIC models | | Proposed models | | |
|---|---|---|---|---|---|
| JPEG2000 | [8] | [15] | $\alpha = 0.1$ | $\alpha = 0.01$ | $\alpha = 0$ |
| opt-MSE | -17.48 | -14.48 | -16.11 | -12.72 | -13.24 |
| opt-MS-SSIM | -62.48 | -64.88 | -63.30 | -62.51 | -63.73 |

Conventional scalable codecs cost 15-25% overhead bits to add a scalability layer in comparison with equivalent non-scalable codecs [26]. Compared to [8], which optimized for MSE, our two-level quality-scalable model with $\alpha = 0.1$ only increases bitrate usage by 1.3%. Furthermore, for $\alpha = 0.1$, our codec outperforms our reimplementation of the non-scalable OctConv model in [15]. For the second row comparing models optimized for MS-SSIM, the OctConv-based LIC models show coding gains compared to [8], which uses only regular convolution. We attribute these gains to the structural ability of OctConv in efficiently capturing the spatial structure of their input in a manner reminiscent of Wavelet transforms.

### B. Quality scalability and frequency analysis

Fig. 6 shows a visual example of quality scalablility in action, where an input image $\mathbf{X}$ is compressed using our MSE-optimized FLIC model trained with small $\lambda = 0.01$ and $\alpha = 0.01$ to make the quality degradation more apparent. Each of the columns shows an image, its Fourier domain, and two select enlarged ROIs. The image in each column in Fig. 6 is: (a) $\mathbf{X}$, the original uncompressed image, (b) $\widehat{\mathbf{X}}$, reconstructed using both base and enhancement bitstreams, (c) $\widehat{\mathbf{X}}_{\text{base}}$, reconstructed using only the base bitstream, and

(d) $\widehat{\mathbf{X}}_{\text{ROI}}$, reconstructed using the base bitstream along with the enhancement bitstream containing only selected ROIs. Since the base bitstream primarily contains low-frequency information, much of the high-frequency information is missing from $\widehat{\mathbf{X}}_{\text{base}}$. Indeed, the Fourier domain for $\widehat{\mathbf{X}}_{\text{base}}$ shows much less energy in the high-frequency spectrum compared to $\mathbf{X}$ and $\widehat{\mathbf{X}}$.

To demonstrate our model's capability for ROI-based quality scalability, we select two ROIs for enhancement: the man in the center of the boat, and the text on the boat. These regions are blocky and unreadable in $\widehat{\mathbf{X}}_{\text{base}}$. By including these regions within the enhancement bitstream, the decoder effectively reconstructs these regions within $\widehat{\mathbf{X}}_{\text{ROI}}$, with a quality equivalent to the same regions in $\widehat{\mathbf{X}}$ by PSNR.

### V. CONCLUSION

We presented a novel frequency-aware learned image compression (FLIC) framework that uses our newly introduced WeOctConv layer. The WeOctConv layer is designed to optimize the separation of spatial frequencies into two latent representations, enabling our FLIC to serve two-level quality-scalable coding with minimal overhead bits. Furthermore, our method efficiently balances the amount of information between the low- and high-frequency latent representations using a scale factor during training. This allows it to achieve flexible RD trade-offs for the base bitstream while having minimal impact on overall coding gain. Finally, we demonstrated the potential of our approach in the context of ROI-based quality enhancement by utilizing partial information from the enhancement latent representation, without requiring any extra retraining.

## REFERENCES

[1] K. Sayood, *Introduction to data compression*, Morgan Kaufmann, 2017.

[2] J. M. Shapiro, "Embedded image coding using zerotrees of wavelet coefficients," *IEEE Trans. Signal Process.*, vol. 41, no. 12, pp. 3445–3462, 1993.

[3] C. Christopoulos, A. Skodras, and T. Ebrahimi, "The jpeg2000 still image coding system: an overview," *IEEE Trans. Consum. Electron.*, vol. 46, no. 4, pp. 1103–1127, 2000.

[4] H. Schwarz, D. Marpe, and T. Wiegand, "Overview of the scalable video coding extension of the h. 264/avc standard," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 17, no. 9, pp. 1103–1120, 2007.

[5] J. M. Boyce, Y. Ye, J. Chen, and A. K. Ramasubramonian, "Overview of shvc: Scalable extensions of the high efficiency video coding standard," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 26, no. 1, pp. 20–34, 2015.

[6] J. Ballé, V. Laparra, and E. P. Simoncelli, "Density modeling of images using a generalized normalization transformation," in *Proc. ICLR*, 2016.

[7] G. Ian, B. Yoshua, and C. Aaron, *Deep Learning*, MIT Press, 2016.

[8] J. Ballé, V. Laparra, and E. P. Simoncelli, "End-to-end optimized image compression," in *Proc. ICLR*, 2017.

[9] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," in *Proc. ICLR*, 2018.

[10] D. Minnen, J. Ballé, and G. D. Toderici, "Joint autoregressive and hierarchical priors for learned image compression," *Adv. Neural Inf. Process. Syst.*, vol. 31, 2018.

[11] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, "Learned image compression with discretized gaussian mixture likelihoods and attention modules," in *Proc. IEEE CVPR*, 2020.

[12] Int. Telecommun. Union-Telecommun. (ITU-T) and Int. Standards Org./Int/Electrotech. Commun. (ISO/IEC JTC 1), "High efficiency video coding," Rec. ITU-T H.265 and ISO/IEC 23008-2, 2019.

[13] Int. Telecommun. Union-Telecommun. (ITU-T) and Int. Standards Org./Int/Electrotech. Commun. (ISO/IEC JTC 1), "Versatile video coding," Rec. ITU-T H.266 and ISO/IEC 23090-3, 2020.

[14] Y. Xie, K. L. Cheng, and Q. Chen, "Enhanced invertible encoding for learned image compression," in *Proc. ACM Int. Conf. Multimed.*, 2021.

[15] M. Akbari, J. Liang, J. Han, and C. Tu, "Learned bi-resolution image coding using generalized octave convolutions," in *Proc. AAAI*, 2021, pp. 6592–6599.

[16] Y. Chen, H. Fan, B. Xu, Z. Yan, Y. Kalantidis, M. Rohrbach, S. Yan, and J. Feng, "Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution," in *Proc. IEEE ICCV*, 2019, pp. 3435–3444.

[17] M. Wolter and J. Garcke, "Adaptive wavelet pooling for convolutional neural networks," in *Int. Conf. Artif. Intell. Stat.* PMLR, 2021, pp. 1936–1944.

[18] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proc. IEEE CVPR*, 2016, pp. 1874–1883.

[19] J. Bégaint, F. Racapé, S. Feltman, and A. Pushparaja, "Compressai: a pytorch library and evaluation platform for end-to-end compression research," *arXiv preprint arXiv:2011.03029*, 2020.

[20] E. Kodak, "Kodak lossless true color image suite (PhotoCD PCD0992)," http://r0k.us/graphics/kodak.

[21] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *Proc. IEEE Asilomar Conf. Signals, Systems & Computers*, 2003, vol. 2, pp. 1398–1402.

[22] T. Xue, B. Chen, J. Wu, D. Wei, and W. T. Freeman, "Video enhancement with task-oriented flow," *Int. J. Comput. Vision*, vol. 127, no. 8, pp. 1106–1125, 2019.

[23] S. Tomar, "Converting video formats with ffmpeg," *Linux Journal*, vol. 2006, no. 146, pp. 10, 2006.

[24] "JPEG2000 reference software," [Online]: https://github.com/uclouvain/openjpeg, Accessed: 2022-05-04.

[25] G. Bjøntegaard, "VCEG-M33: Calculation of average PSNR differences between RD curves," in *Video Coding Experts Group (VCEG)*, Apr. 2001.

[26] J. M. Boyce, Y. Ye, J. Chen, and A. K. Ramasubramonian, "Overview of SHVC: scalable extensions of the high efficiency video coding standard," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 26, no. 1, pp. 20–34, 2016.