

LC-FDNet: Learned Lossless Image Compression with Frequency Decomposition Network

Hochang Rhee¹, Yeong Il Jang¹, Seyun Kim², Nam Ik Cho¹

¹Department of ECE, INMC, Seoul National University, Seoul, Korea

²Gauss Labs Inc.

hochang, jyicu@ispl.snu.ac.kr, seyun.kim@gausslabs.ai, nicho@snu.ac.kr

Abstract

Recent learning-based lossless image compression methods encode an image in the unit of subimages and achieve comparable performances to conventional non-learning algorithms. However, these methods do not consider the performance drop in the high-frequency region, giving equal consideration to the low and high-frequency areas. In this paper, we propose a new lossless image compression method that proceeds the encoding in a coarse-to-fine manner to separate and process low and high-frequency regions differently. We initially compress the low-frequency components and then use them as additional input for encoding the remaining high-frequency region. The low-frequency components act as a strong prior in this case, which leads to improved estimation in the high-frequency area. In addition, we design the frequency decomposition process to be adaptive to color channel, spatial location, and image characteristics. As a result, our method derives an image-specific optimal ratio of low/high-frequency components. Experiments show that the proposed method achieves state-of-the-art performance for benchmark high-resolution datasets.

1. Introduction

As the need for high-quality images is increasing, the importance of image compression is growing accordingly. Driven by the development of deep neural networks (DNNs), there has been remarkable progress in computer vision and image processing, including lossy [3, 4, 8, 10–17, 21, 23–25, 27, 28, 30, 32–34, 40, 44–48, 53, 54] and lossless image compression [3, 9, 18, 19, 29, 31, 35, 38, 39, 41, 42, 49]. Although lossy compression is generally preferred, lossless compression is also necessary for many applications. Lossless compression is especially required for medical images, scientific images, technical drawings, and artistic photos. While methods such as JPEG2000 (lossless mode) [37] employ transform coding with discrete wavelet transform

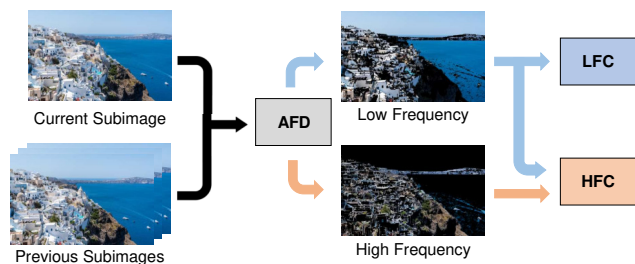


Figure 1. Our LC-FDNet consists of Adaptive Frequency Decomposition (AFD), Low-Frequency Compressor (LFC), and High-Frequency Compressor (HFC). The current subimage is split into low/high-frequency regions through the AFD. The LFC first compresses the low-frequency region, and then the HFC compresses the high-frequency area using the low-frequency as strong prior.

(DWT), most of the standard/non-standard lossless compression methods [2, 5, 7, 51] use predictive coding. The standard predictive coding scheme uses a closed-loop prediction where the current pixel is estimated and compressed using the previously encoded samples.

In this sense, early learning-based lossless compression algorithms [29, 31, 35, 38, 41, 49] design DNNs as autoregressive models. They rely on the strong power of DNNs in estimating the probability distribution of a pixel conditioned on the previous samples. For example, PixelRNN [49], PixelCNN [35], and PixelCNN++ [41] compress each pixel sequentially, where the probability distribution is predicted conditioned on all previous pixels. However, these methods require neural network computations for the number of whole pixels, leading to an impractical inference time.

To achieve practicality, recent works [29, 31, 38] process the encoding in the unit of an entire image or subimages rather than individual pixels. These methods derive the probability distribution of a subimage conditioned on the previously encoded subimages, or the distribution of a whole image conditioned on the lossy compressed image. They show reduced and practical computation

time compared to pixel-wise encoding methods. However, these methods consider the low and high-frequency regions equally, giving the same encoding strategies to the regions of different characteristics. In general, it is difficult to obtain optimal performance in high-frequency regions near an edge or texture where the pixel values change rapidly.

We address this challenge and propose Lossless Compression with Frequency Decomposition Network (LC-FDNet) illustrated in Fig. 1, which consists of Adaptive Frequency Decomposition (AFD), Low-Frequency Compressor (LFC), and High-Frequency Compressor (HFC). We also decompose an image into subimages based on our unique decomposition scheme, where the first subimage is compressed by a conventional lossless compressor. Then, the rest subimages are sequentially compressed by Fig. 1. Using the previously encoded and current subimages as the input, the AFD decomposes the input subimage into low and high-frequency regions, and the compressors (LFC and HFC) encode low and high-frequency regions differently. Since the low-frequency region is typically well predicted, we first compress the low-frequency components. On the other hand, high-frequency regions usually exhibit relatively large prediction errors, and hence we encode them separately with additional priors, which are the encoded low-frequency pixels. That is, we feed the low-frequency components as additional input for compressing the high-frequency region.

For the image-specific frequency decomposition, the AFD generates *error variance map* and *error variance thresholds*. Error variance map can be comprehended as the magnitude of the prediction error produced by the network. By thresholding the error variance map with the error variance threshold, we can classify the pixels into low and high-frequency ones. Since the error variance differs depending on the channel, spatial location, and image characteristics, we design the threshold to be adaptive to those factors. This drives the frequency decomposition process to be image-specific, where different threshold values are derived depending on the image property. Experiments show that the proposed method achieves state-of-the-art performance for benchmark high-resolution datasets with reasonable inference time.

In summary, the main contributions are as follows:

- We propose a lossless image compression framework that compresses in a coarse-to-fine manner, using the low frequency components to boost the performance in high-frequency regions.
- We design the frequency decomposition process to be adaptive to channel, spatial location, and image characteristics. Hence, the encoding becomes image-specific, improving the compression performance.

- Our method achieves state-of-the-art performance for benchmark high-resolution datasets with reasonable inference time.

2. Related Works

Pixel-wise Lossless Compression Learning-based lossless compression methods generally adopt an autoregressive model. Early methods proceeded the encoding in the pixel unit, where each pixel is compressed based on the previously encoded ones. For example, PixelRNN [49] and PixelCNN [35] modeled a pixel as the product of conditional distributions $p(\mathbf{x}) = \prod_i p(x_i|x_1, \dots, x_{i-1})$, where x_i is a single pixel. PixelCNN++ [41] was proposed as an advancement of the above works and achieved performance enhancement along with faster time. They modeled the pixels as a discretized logistic mixture likelihood, used down-sampling to capture structure at multiple resolutions, and introduced additional short-cut connections. Despite these factors, PixelCNN++ still maintains the inherent limitation of autoregressive models, *i.e.*, network computation is required for each pixel, requiring impractical inference time.

Subimage-wise Lossless Compression For the lossless compression in a reasonable time, recent works perform the encoding in the unit of an entire image or subimages. Each of these methods has its unique strategy for converting an image into subimages. MS-PixelCNN [38] first proposed a parallelized PixelCNN using a hierarchical encoding scheme. Specifically, the input image is explicitly divided into four subimages depending on the spatial location, and the distribution of a subimage is conditioned on the previously encoded subimages. However, they used PixelCNNs for modeling the dependency between the subimages, which still required impractical time. L3C [29] proposed a practical compression framework that utilizes a hierarchical probabilistic model. The subimages are implicitly modeled by a neural network and each subimage is conditioned on the subimage of the previous scale. Here, the initial subimage is assumed as a uniform distribution. RC [31] can be seen as a method that divides the image into two parts: lossy compressed image and its residuals. The probability distribution of the residuals is modeled based on lossy compression.

3. Method

3.1. Overview

The overall procedure of our method is illustrated in Fig. 2. Given the input image $x \in \mathcal{R}^{H \times W \times 3}$, we first convert the RGB image into a YUV format through a reversible color transform [36]. Then we split the image in a channel-wise and spatial-wise manner. Specifically, we divide the input image into 12 subimages $x_{c,s} \in \mathcal{R}^{\frac{H}{2} \times \frac{W}{2} \times 1}$, where c

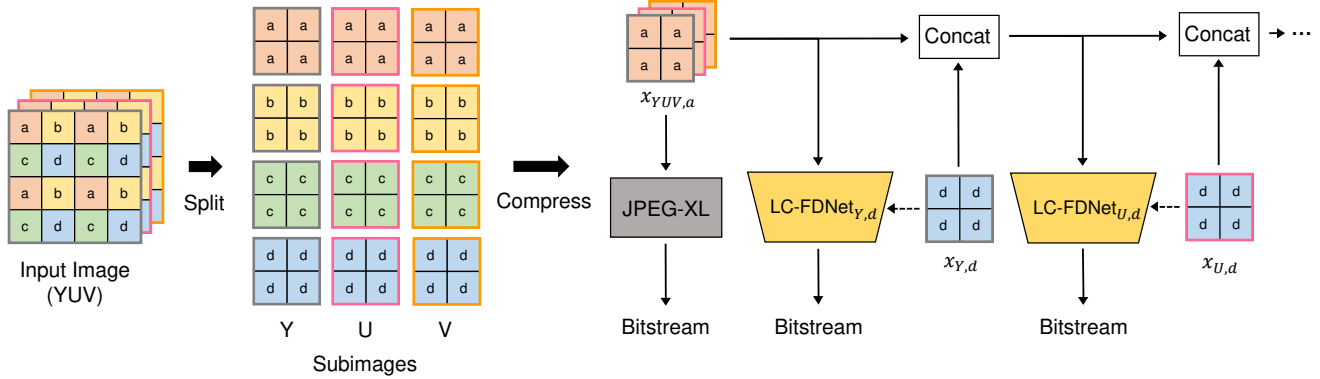


Figure 2. The framework of our compression scheme. Depending on the spatial location, each pixel is grouped as either a, b, c, d . The input image is split into subimages, which are sequentially compressed. The subimage $x_{YUV,a}$ is initially encoded using a conventional compression algorithm. The remaining subimages are compressed through deep networks, which receive the previously encoded subimages as input and compress the current subimage. The dotted arrow denotes that the corresponding subimage is currently being compressed. The compressed subimage is then used as an additional input for encoding the next subimage.

denotes the channel index ($c \in \{Y, U, V\}$) and s denotes the spatial location index ($s \in \{a, b, c, d\}$). The subimages $x_{YUV,a} = \{x_{Y,a}, x_{U,a}, x_{V,a}\}$ are first compressed using a conventional compression algorithm. Then, the remaining subimages are compressed one by one with our LC-FDNet, where previously encoded subimages are used as input. The order of the subimages to be compressed will be explained in Sec.3.3.

3.2. Reversible Color Transform

In general, RGB images have significant correlations between the color channels. Most standard image/video compression methods adopt YUV transformation to decorrelate the color channels and enhance the compression efficiency. In the case of lossless compression, the YUV transformation must be itself lossless, where the inverse of the YUV back to the RGB should be lossless in integer arithmetic. In this paper, we adopt the reversible color transform proposed in [36] since it well approximates the conventional YUV transformation. Note that the Y channel is expressed in 8 bits and UV channels are expressed in 9 bits.

3.3. Framework

After the reversible color transformation, the input image is divided into subimages depending on the color channel and the spatial location. Fig. 2 shows how we categorize the pixels into four groups (a, b, c, d) depending on the spatial location. Pixels in the odd row and odd column are categorized as a , odd row and even column as b , and so on.

We compress 12 subimages in total, where the compression of each subimage is conditioned on the previously encoded subimages. To be specific, for the compression of the N -th subimage $y \in \mathcal{R}^{\frac{H}{2} \times \frac{W}{2} \times 1}$, we concatenate the already encoded $N - 1$ subimages and use it as the input, which we

denote as $x_{in} \in \mathcal{R}^{\frac{H}{2} \times \frac{W}{2} \times (N-1)}$. We neglect the notation of the subimage index N for the sake of simplicity. In this scenario, the order of the subimages is critical to compression efficiency. The compression performance is improved as the correlation among the input and the N -th subimage increases. For instance, encoding subimage $x_{Y,d}$ is much easier when it is conditioned on $x_{Y,a}$ rather than $x_{V,b}$.

We design the order of the subimages considering the following two factors : 1) color channel and 2) spatial location. In terms of the color channel, we arrange the order as $Y \rightarrow U \rightarrow V$. This is a straightforward choice since the Y channel contains more significant features than U and V . Considering the spatial location, we design the network to proceed in the order of $a \rightarrow d \rightarrow b \rightarrow c$. We figure that this is a better design choice compared to MS-PixelCNN [38], where they progress in the order of $a \rightarrow b \rightarrow c \rightarrow d$. Comparing $a \rightarrow d$ and $a \rightarrow b$, we see that d fully utilizes the information of a both horizontally and vertically. In contrast, acquiring b conditioned on a may have more benefit in terms of the horizontal axis, but lacks to fully utilize the vertical components. In conclusion, we proceed the compression of the subimages in the order of $x_{Y,a} \rightarrow x_{U,a} \rightarrow x_{V,a} \rightarrow x_{Y,d} \rightarrow x_{U,d} \rightarrow \dots \rightarrow x_{V,c}$.

For the compression of the initial subimage $x_{YUV,a}$, we adopt a conventional lossless compression algorithm, similar to RC [31]. Prior works [29, 33] provide the initial prior as uniform distribution or unit Gaussian distribution. Although DNNs estimate conditional probability distributions very well, the strength is limited when weak priors are given. Conventional algorithms instead show competitive performance in this environment, and hence we use them for compressing the initial subimage $x_{YUV,a}$. Specifically, we adopt JPEG-XL [2] which yields state-of-the-art performance among the conventional compression algorithms.

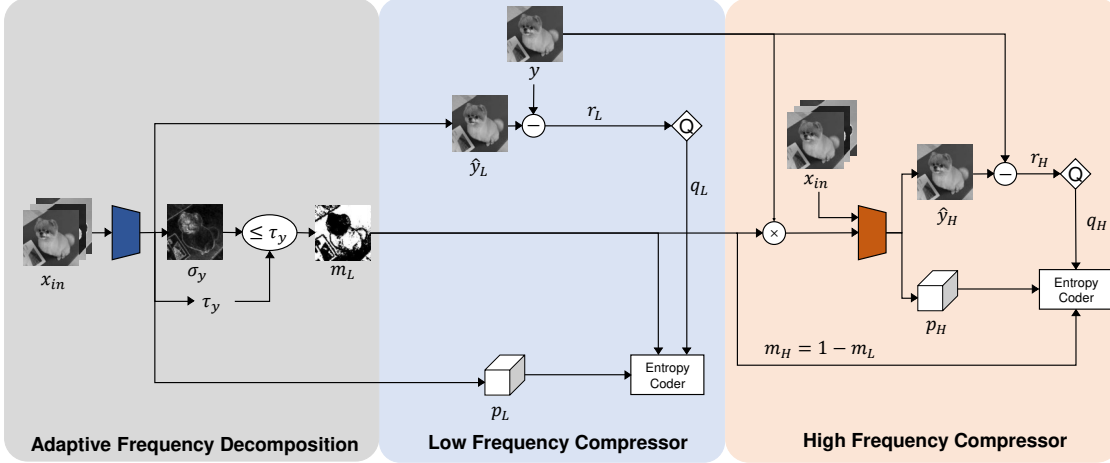


Figure 3. The architecture of LC-FDNet. In this figure, we consider the case of compressing $y = x_{Y,d}$ given $x_{in} = x_{YUV,a}$. AFD part first receives x_{in} and determines each pixel as belonging to either low or high-frequency regions, using error variance map σ_y and error variance threshold τ_y . Afterward, LFC encodes the low-frequency region of subimage y . HFC then receives the encoded low-frequency region as additional input and compresses the remaining high-frequency region. The decoding process is provided in the Supplementary.

3.4. Architecture

In this section, we present the architecture of LC-FDNet illustrated in Fig. 3, which shows the details of Fig. 1. The goal is to compress the N -th subimage y given the input x_{in} , which is the concatenation of $N - 1$ previously encoded subimages. Note that LC-FDNet is required for each of the subimages, resulting in 9 LC-FDNet in total (since 3 subimages are encoded by JPEG-XL). These networks do not share the parameters since each of them is specific for each subimage.

Throughout the paper, notations L and H denote low and high-frequency, respectively. We first explain the notations in the AFD and LFC parts.

Subimage Prediction $\hat{y}_L \in \mathcal{R}^{\frac{H}{2} \times \frac{W}{2} \times 1}$ is the network prediction of y , where better prediction yields more compact compression. Given the prediction, the residual is computed as $r_L = \hat{y}_L - y$, which is the difference between the ground truth subimage and the prediction. Since the obtained residual is not in the form of integers, we quantize (round) the residual. It is denoted as q_L , which is then passed to the entropy coder.

Probability Distribution p_L is the estimated probability distribution of the quantized residual q_L . We directly estimate the probability distribution as the probability mass function (pmf). Hence, the dimension of p_L is $\frac{H}{2} \times \frac{W}{2} \times C$, where C is 511 for the Y channel, and 1021 for the U, V channel. Softmax operation is applied before deriving p_L , so that the probabilities sum up to 1.

Error Variance Map The error variance map $\sigma_y \in \mathcal{R}^{\frac{H}{2} \times \frac{W}{2} \times 1}$ represents the estimation of the prediction error

magnitude generated by the network. We design the error variance map to follow the magnitude of the prediction error through the following loss:

$$L_{ev} = \|\sigma_y - |y - \hat{y}_L|\|_1. \quad (1)$$

Here, each value in the map can be interpreted as the prediction error variance at the corresponding pixel. A large value implies that the network is likely to make a large prediction error at the point, which means that the pixel belongs to a high-frequency region. Similarly, smooth regions *i.e.*, low-frequency regions yield low error variance values.

Error Variance Threshold With the obtained error variance map, we apply a simple thresholding to categorize each pixel into two classes; pixels in low or high-frequency regions. However, the threshold value should be adjusted depending on the channel, spatial location and image characteristics. For instance, the error variance is typically larger in the Y channel compared to U and V . Thus the threshold should be larger in the Y . Therefore, instead of a fixed threshold, we let the network derive a specific error variance threshold $\tau_y \in \mathcal{R}$ for each subimage. Note that 9 threshold values are derived for a single input image. With σ_y and τ_y in hand, we acquire the low-frequency mask as

$$m_L^i = \begin{cases} 1 & \text{if } \sigma_y^i \leq \tau_y \\ 0 & \text{else,} \end{cases} \quad (2)$$

where i denotes the pixel index. m_L serves as an indicator of which components are considered in the low-frequency region.

The quantized residual q_L , corresponding probability distribution p_L , and the low-frequency mask m_L are passed

to the entropy coder. We compress only the low-frequency components *i.e.*, pixels corresponding to $m_L^i = 1$. It can be assumed that pixels belonging to low-frequency regions will have marginal performance enhancement even when additional information is given, specifically when it is compressed in the HFC. Instead, the compression efficiency gain is significant when these components serve as the additional input.

After the compression of low-frequency regions in LFC, HFC encodes the remaining high-frequency regions. Besides x_{in} , HFC additionally receives the low-frequency component of the currently encoding subimage $y \odot m_L$. From the input, HFC generates the following two outputs: 1) \hat{y}_H : the prediction of y , 2) p_H : the probability distribution of the quantized residual q_H . Since low-frequency components serve as a strong prior for the high-frequency components, HFC can make more precise predictions. In addition, the variance of the probability distribution is reduced, leading to compression efficiency.

The pipeline of HFC is similar to that of LFC. The quantized residual q_H , corresponding probability distribution p_H , and the high-frequency mask $m_H = 1 - m_L$ are fed to the entropy coder. Note that HFC can ignore the estimation of low-frequency components and only focus on the high-frequency ones.

3.5. Loss Function

LC-FDNet is trained with the following three losses: 1) Error variance loss defined as Eq. 1, 2) reconstruction loss, and 3) bitrate loss.

Reconstruction Loss We define reconstruction loss as the L1 loss between the ground truth and the predicted subimage:

$$L_{rec} = m_L \cdot \|y - \hat{y}_L\|_1 + m_H \cdot \|y - \hat{y}_H\|_1. \quad (3)$$

Note that we multiply the corresponding frequency mask to the prediction error of LFC and HFC. This lets only the low-frequency components contribute to the reconstruction loss of LFC, and similarly for HFC. This makes the LFC/HFC to be specified for low/high-frequency regions, respectively. Although the reconstruction loss is often neglected in other researches, we figure that adopting this loss leads to stable training and performance enhancement.

Bitrate Loss Bitrate loss is used to minimize the cross-entropy between the real probability distribution of the quantized residual (p_{q_L}, p_{q_H}) and the estimated (p_L, p_H), respectively. Formally, it is defined as

$$L_{br} = m_L \cdot \|\log p_L(q_L)\|_1 + m_H \cdot \|\log p_H(q_H)\|_1. \quad (4)$$

The probability distributions p_L and p_H are trained to classify the corresponding quantized residuals (symbols) q_L and

q_H by the cross-entropy loss. This is equivalent to the expected bits per symbol and thus we can directly minimize the coding cost. To restrict the contribution of each frequency component as in the reconstruction loss, we multiply the frequency masks to the corresponding probability distribution.

Altogether, we train our network with the loss:

$$L = L_{rec} + \lambda_{br} L_{br} + \lambda_{ev} L_{ev} \quad (5)$$

where λ_{ev} and λ_{br} are the balancing hyperparameters. In our experiments, we set both λ_{ev} and λ_{br} as 1.

4. Experiments

4.1. Experimental Setup

Implementation Detail The detail of the network architecture is provided in the Supplementary Material. For the quantization, we use the round function *i.e.*, $q = \text{round}(r)$. The derivative is zero except at integers, which cannot be used in gradient-based optimization. Therefore, we approximate the round function as simple STE [6] *i.e.*, $q = r$ in the backward pass since [10] has shown that different quantization approximation methods have a minor effect on the compression performance. The same problem is introduced when deriving m_L with Eq. 2. This is approximated as $m_L = \text{sigmoid}(-(\sigma_y - \tau_y))$ in the backward pass. For our entropy coder, we use “torchac,” which is a fast arithmetic coding library for PyTorch developed by the authors of L3C [29].

Dataset We validate our method on three benchmark datasets, CLIC.m, CLIC.p and DIV2K. CLIC mobile (CLIC.m) and CLIC professional (CLIC.p) are datasets released as part of the “Workshop and Challenge on Learned Image Compression” [52]. CLIC.m consists of 61 evaluation images which are taken using mobile phones. CLIC.p contains 41 evaluation images which are taken by DSLRs. Most of the images in the CLIC datasets are 2K resolution, but some of them are low resolution as far as 512×384 . DIV2K [1] is a super-resolution dataset that consists of 2K resolution high-quality images, where 100 images are provided for evaluation.

Training We train our network with Flickr2k [26] dataset that provides 2,000 high-quality images for training. We randomly extract a patch of size 128×128 from the input image during training. Adam optimizer [20] is used for the training, with a batch size of 24 for 3,000 epochs. The learning rate is initially set as 1×10^{-3} and decays by a factor of 0.1 every 1,000 epochs. The training takes 36 hours when trained on a GeForce GTX 1080 Ti.

Evaluation We compare our method for both learned and non-learned compression algorithms. We compare against

Table 1. Comparison of our method with other non-learning and learning-based codecs on high-resolution benchmark dataset. We measure the performances in bits per pixel (bpp). Best performance is highlighted in bold and the second-best performance is denoted with *. The difference in percentage to our method is highlighted in green.

Method	CLIC.m	CLIC.p	DIV2K
PNG [7]	11.79 +70.9%	11.79 +50.0%	12.69 +55.9%
JPEG-LS [51]	7.59 +10.0%	8.46 +7.6%	8.97 +10.2%
JPEG2000 [37]	8.13 +17.8%	8.79 +11.8%	9.36 +15.0%
WebP [50]	8.19 +17.8%	8.70 +11.8%	9.33 +15.0%
BPG [5]	8.52 +23.5%	9.24 +17.6%	9.84 +20.9%
LCIC [19]	7.88 +14.2%	9.02 +14.8%	9.35 +14.9%
FLIF [43]	7.44 +7.8%	8.16 +3.8%	8.73 +7.2%
JPEG-XL [2]	7.20* +4.3%	8.19 +4.2%	8.49 +4.3%
L3C [29]	7.92 +14.8%	8.82 +12.2%	9.27 +13.9%
RC [31]	7.62 +10.4%	8.79 +11.8%	9.24 +13.5%
Near-Lossless [3]	7.53 +9.1%	7.98* +1.5%	8.43* +3.6%
Ours	6.90	7.86	8.14

the following conventional lossless image codecs: PNG [7], JPEG-LS [51], JPEG2000 [37], WebP [50], BPG [5], LCIC [19], FLIF [43] and JPEG-XL [2]. As for the learned methods, we consider L3C [29], RC [31], and Near-Lossless [3]. L3C and RC are trained with Open Images dataset [22] consisting of 300,000 images. Near-Lossless is trained with the DIV2K dataset. We use *bits per pixel* (bpp) as the evaluation metric, where lower bpp indicates better compression performance.

4.2. Compression Result

Table 1 presents the comparisons on the described evaluation sets. It can be seen that our method shows superior performance to both engineered and learning-based codecs. Considering DIV2K, our method achieves a 3.6% gain compared to Near-Lossless. In the case of CLIC.m, non-learning codecs such as FLIF and JPEG-XL outperform existing learning-based methods. Thus, it can be interpreted that learning-based methods are difficult to be generalized to CLIC.m. Nevertheless, our method achieves state-of-the-art performance and outperforms JPEG-XL by 4.3%. Finally, for CLIC.p, our method shows the best performance achieving 1.5% gain against Near-Lossless.

In Table 2, we report the compression result for each subimage, *i.e.*, each channel and spatial location. Considering the spatial location, we observe that the compression efficiency enhances in the order of $a \rightarrow d \rightarrow b \rightarrow c$. This is straightforward since more information is supplied as we proceed in the above order. From the perspective of channels, better compression is presented in the order of $Y \rightarrow UV$. This is due to the color transform that reduces the variance in UV channels. Moreover, V channel shows a slight improvement compared to U since we use additional input U when encoding the V .

Table 2. Compression result of each subimage for the DIV2K dataset. Compression performance of subimage $x_{YUV,a}$ is the result of JPEG-XL.

bpp	a	d	b	c
Y	-	0.95	0.81	0.79
U	-	0.58	0.45	0.44
V	-	0.57	0.44	0.43
Total	2.68	2.10	1.70	1.66

4.3. Inference Time

We measure the inference time required for encoding a 512×512 image on an GeForce GTX 1080 Ti. The decoding time is the same as the encoding time for our method. First, the compression of the initial subimage using JPEG-XL requires 199 ms. The forward pass for achieving the quantized residual, probability distribution, and frequency mask takes 33 ms. Finally, the arithmetic coding using torchac requires 611 ms. In total, our method requires 841 ms. Note that 72% of the time is consumed in arithmetic coding, which can be shortened if better PyTorch-friendly entropy coder is developed in the future.

4.4. Adaptive Frequency Decomposition Analysis

We quantitatively and qualitatively demonstrate that the error variance threshold is adaptive to the channel, spatial location, and image characteristics. We first show that the error variance threshold is adaptive to the channel and spatial location through Table 3. In general, the threshold value decreases in the order of $d \rightarrow b \rightarrow c$ and $Y \rightarrow UV$. This is consistent with the order of compression efficiency. If the subimage is more predictable, the overall values in the variance map tend to decrease. In this case, the error vari-

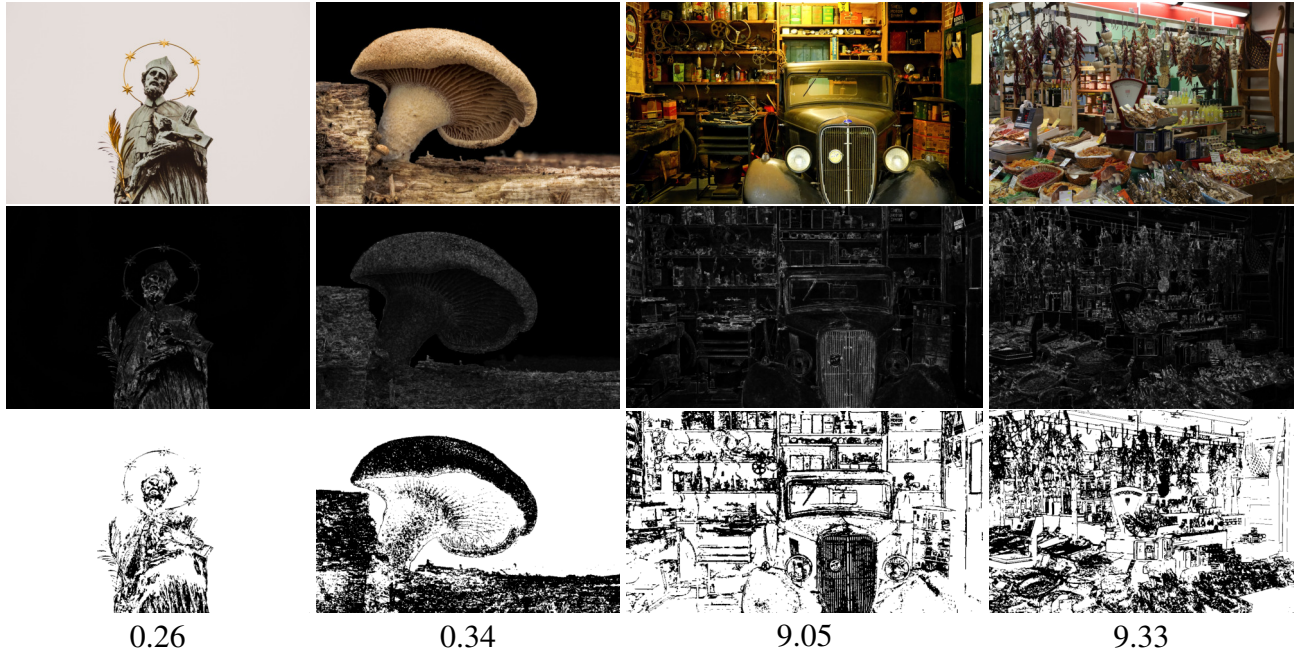


Figure 4. From top to down are the visualizations of an input image, error variance map, low-frequency mask, and error variance threshold. These elements are visualized for the case of the Y channel and d location. We choose the samples from the DIV2K dataset that have the smallest and largest τ_y . The error variance map is magnified by 5 for visualization.

Table 3. Error variance threshold value for each subimage. Since the threshold is image-specific, we average the threshold for all the images in DIV2K dataset.

τ_y	d	b	c
Y	3.57	2.84	2.72
U	2.68	2.24	2.21
V	2.66	2.30	2.24

ance threshold should also decrease to balance the low to high-frequency ratio.

Next, we show that the error variance threshold is adaptive to the image characteristics. Fig. 4 shows the outputs generated from LFC. The first two samples contain a large portion of the smooth background and a single object. These samples produce a small threshold value of 0.26 and 0.34. In contrast, the last two samples are more complicated than the preceding ones and introduce many high-frequency components. These samples generate a large threshold value of 9.05 and 9.33. We interpret that the error variance threshold is proportional to the number of high-frequency components an image contains.

We also verify the above conclusion quantitatively. We figure that images with many high-frequency components tend to introduce large values in the error variance map. In addition, these images result in low compression rate. Hence, for samples in DIV2K, we plot the error variance

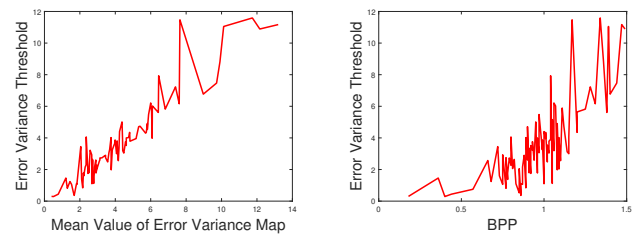


Figure 5. Graph of image characteristic versus error variance threshold for DIV2K dataset. We use the samples of DIV2K, channel of Y , and spatial location of d .

threshold against the mean value of the error variance map and bpp in Fig. 5. It can be observed that the error variance threshold and the two components have a positive correlation. In conclusion, the error variance threshold is adaptive to image characteristics, where the threshold value increases as more high-frequency components are present.

4.5. Ablation Study

Several ablation experiments are performed to analyze each component of LC-FDNet. We demonstrate the contribution of each component in Table 4 by excluding the components one by one. The comparing networks are evaluated on the DIV2K dataset. We exclude the portion of JPEG-XL (2.68 bpp) in computing the compression performance.

Table 4. Ablation study of our method on DIV2K dataset. C2F refers to the coarse-to-fine network. \checkmark indicates that the corresponding element is used.

C2F	Adaptive τ	Loss Masking	bpp
	\checkmark		5.78 +5.8%
\checkmark		\checkmark	5.61 +2.7%
\checkmark	\checkmark		5.57 +2.0%
\checkmark	\checkmark	\checkmark	5.46

Coarse to Fine We first demonstrate the effect of proceeding in a coarse-to-fine manner. We design a comparing network that compresses the low and high-frequency components together so that the high-frequency components do not benefit from the low ones. Hence, the network only outputs the subimage prediction and probability distribution. We match the number of parameters for both networks to demonstrate that the performance gain does not come from the difference in the network size. The result (first row of Table 4) shows that we can have a 5.8% performance gain by the coarse-to-fine processing. Hence, we can conclude that low-frequency components act as strong priors for estimating the high-frequency components.

Adaptive Error Variance Threshold We show that letting the error variance threshold be adaptive to image characteristics leads to performance enhancement. Specifically, we train a network with fixed τ for every subimage. Since our framework is sensitive to the value of τ , we should carefully set the threshold value for a fair comparison. Hence, we use the average τ_y of the DIV2K validation set derived from our full model as our fixed τ . The second row of the table shows that the performance decreases by 2.7% with the fixed τ .

Loss Masking We verify that multiplying the corresponding frequency mask in Eq. 3 and Eq. 4 has a valid contribution. Here, we train a network without the multiplication of frequency mask. In other words, the LFC and HFC share the same objective and are not frequency-specific. In this scenario, a performance drop of 2.0% is observed as in the third row, indicating that assigning frequency-specific roles to LFC and HFC has a positive influence.

4.6. Frequency Component Analysis

We show the performance enhancement in low and high-frequency regions separately in Table 5 for further analysis of our system. We first compare our full model (C2F) against the model without the coarse-to-fine processing (w/o C2F) as in the ablation study. Compared to the network that proceeds without the coarse-to-fine processing, the low-frequency components have a 3.4% performance gain, whereas high-frequency components have an 11.4% increase. This implies that the high-frequency components benefit significantly from the coarse-to-fine processing. The

Table 5. Performance gain on low and high-frequency regions. F2C refers to the network proceeding in a fine-to-coarse manner.

Method	Low-Freq	High-Freq	Total
w/o C2F	3.93 +3.4%	1.85 +11.4%	5.78 +5.8%
F2C	3.78 -0.5%	1.84 +10.8%	5.62 +2.9%
C2F	3.80	1.66	5.46

low-frequency components indeed act as strong conditions for the estimation of high-frequency ones, as intended.

We train an additional network that proceeds in a fine-to-coarse manner (F2C). That is, we compress the high-frequency components first and utilize them for encoding the low-frequency components. From Table 5, we observe that the performance gain is 0.5% in the low-frequency area, which is minor. In contrast, there is a considerable performance drop of 10.8% in the high-frequency components. Altogether, there is a total performance drop of 2.9% when proceeding in a fine-to-coarse manner. Although low-frequency components take up a large portion of an image, the gain is too small to have enough contribution to the overall gain. Thus, we conclude that the design choice of coarse-to-fine manner is indeed favorable.

5. Conclusion

We have proposed LC-FDNet, a lossless image compression framework that decomposes an image into low and high-frequency regions to proceed in a coarse-to-fine manner. We resolved the performance drop in the high-frequency areas by first compressing the low-frequency components and using them as a strong prior for encoding the remaining high-frequency components. Furthermore, we designed the frequency decomposition method to be adaptive to color channel, spatial location, and image characteristics to derive the image-specific optimal ratio of low/high-frequency components. Experiments show that our method achieves state-of-the-art performance for high-resolution datasets. Our code is available at <https://github.com/myideaisgood/LC-FDNet>.

Acknowledgments This work was supported in part by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (2021R1 A2C2007220), in part by the Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT)[NO.2021-0-01343, Artificial Intelligence Graduate School Program (Seoul National University)], and in part by BK21 FOUR program of the Education and Research Program for Future ICT Pioneers, Seoul National University in 2022.

References

- [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. **5**
- [2] Jyrki Alakuijala, Ruud van Asseldonk, Sami Boukourt, Martin Bruse, Iulia-Maria Comşa, Moritz Firsching, Thomas Fischbacher, Evgenii Kliuchnikov, Sebastian Gomez, Robert Obryk, et al. Jpeg xl next-generation image compression architecture and coding tools. In *Applications of Digital Image Processing XLII*, volume 11137, page 111370K. International Society for Optics and Photonics, 2019. **1, 3, 6**
- [3] Yuanchao Bai, Xianming Liu, Wangmeng Zuo, Yaowei Wang, and Xiangyang Ji. Learning scalable ly=constrained near-lossless image compression via joint lossy image and residual compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11946–11955, 2021. **1, 6**
- [4] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016. **1**
- [5] Bellard. Bpg image format. <https://bellard.org/bpg>. **1, 6**
- [6] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. **5**
- [7] Thomas Boutell and T Lane. Png (portable network graphics) specification version 1.0. *Network Working Group*, pages 1–102, 1997. **1, 6**
- [8] Benoit Brummer and Christophe De Vleeschouwer. End-to-end optimized image compression with competition of prior distributions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1890–1894, 2021. **1**
- [9] Sheng Cao, Chao-Yuan Wu, and Philipp Krähenbühl. Lossless image compression through super-resolution. *arXiv preprint arXiv:2004.02872*, 2020. **1**
- [10] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learning image and video compression through spatial-temporal energy compaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10071–10080, 2019. **1, 5**
- [11] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7939–7948, 2020. **1**
- [12] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Variable rate deep image compression with a conditional autoencoder. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3146–3154, 2019. **1**
- [13] Ze Cui, Jing Wang, Shangyin Gao, Tiansheng Guo, Yihui Feng, and Bo Bai. Asymmetric gained deep image compression with continuous rate adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10532–10541, 2021. **1**
- [14] Xin Deng, Wenzhe Yang, Ren Yang, Mai Xu, Enpeng Liu, Qianhan Feng, and Radu Timofte. Deep homography for efficient stereo image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1492–1501, 2021. **1**
- [15] Ge Gao, Pei You, Rong Pan, Shunyuan Han, Yuanyuan Zhang, Yuchao Dai, and Hojae Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021. **1**
- [16] Ge Gao, Pei You, Rong Pan, Shunyuan Han, Yuanyuan Zhang, Yuchao Dai, and Hojae Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021. **1**
- [17] Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. Checkerboard context model for efficient learned image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14771–14780, 2021. **1**
- [18] Emiel Hoogeboom, Jörn WT Peters, Rianne van den Berg, and Max Welling. Integer discrete flows and lossless compression. *arXiv preprint arXiv:1905.07376*, 2019. **1**
- [19] Seyun Kim and Nam Ik Cho. Hierarchical prediction and context adaptive coding for lossless color image compression. *IEEE Transactions on image processing*, 23(1):445–449, 2013. **1, 6**
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. **5**
- [21] Jan P Klopp, Keng-Chi Liu, Liang-Gee Chen, and Shao-Yi Chien. How to exploit the transferability of learned image compression to conventional codecs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16165–16174, 2021. **1**
- [22] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2(3):18, 2017. **6**
- [23] Jooyoung Lee, Seunghyun Cho, and Seung-Kwon Beack. Context-adaptive entropy model for end-to-end optimized image compression. *arXiv preprint arXiv:1809.10452*, 2018. **1**
- [24] Mu Li, Wangmeng Zuo, Shuhang Gu, Jane You, and David Zhang. Learning content-weighted deep image compression. *IEEE transactions on pattern analysis and machine intelligence*, 2020. **1**
- [25] Mu Li, Wangmeng Zuo, Shuhang Gu, Debin Zhao, and David Zhang. Learning convolutional networks for content-weighted image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3214–3223, 2018. **1**
- [26] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single

- image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 5
- [27] Haichuan Ma, Dong Liu, Ning Yan, Houqiang Li, and Feng Wu. End-to-end optimized versatile image compression with wavelet-like transform. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 1
- [28] Ankur Mali, Alexander G Ororbia, and C Lee Giles. The sibling neural estimator: Improving iterative image decoding with gradient communication. In *2020 Data Compression Conference (DCC)*, pages 23–32. IEEE, 2020. 1
- [29] Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Practical full resolution learned lossless image compression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10629–10638, 2019. 1, 2, 3, 5, 6
- [30] Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Conditional probability models for deep image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4394–4402, 2018. 1
- [31] Fabian Mentzer, Luc Van Gool, and Michael Tschannen. Learning better lossless compression using lossy compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6638–6647, 2020. 1, 2, 3, 6
- [32] Fabian Mentzer, George Toderici, Michael Tschannen, and Eirikur Agustsson. High-fidelity generative image compression. *arXiv preprint arXiv:2006.09965*, 2020. 1
- [33] David Minnen, Johannes Ballé, and George Toderici. Joint autoregressive and hierarchical priors for learned image compression. *arXiv preprint arXiv:1809.02736*, 2018. 1, 3
- [34] David Minnen, Johannes Ballé, and George Toderici. Joint autoregressive and hierarchical priors for learned image compression. *arXiv preprint arXiv:1809.02736*, 2018. 1
- [35] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. Conditional image generation with pixelcnn decoders. *arXiv preprint arXiv:1606.05328*, 2016. 1, 2
- [36] Soo-Chang Pei and Jian-Jiun Ding. Improved reversible integer-to-integer color transforms. In *2009 16th IEEE International Conference on Image Processing (ICIP)*, pages 473–476. IEEE, 2009. 2, 3
- [37] Majid Rabbani. Jpeg2000: Image compression fundamentals, standards and practice. *Journal of Electronic Imaging*, 11(2):286, 2002. 1, 6
- [38] Scott Reed, Aäron Oord, Nal Kalchbrenner, Sergio Gómez Colmenarejo, Ziyu Wang, Yutian Chen, Dan Below, and Nando Freitas. Parallel multiscale autoregressive density estimation. In *International Conference on Machine Learning*, pages 2912–2921. PMLR, 2017. 1, 2, 3
- [39] Hochang Rhee, Yeong Il Jang, Seyun Kim, and Nam Ik Cho. Lossless image compression by joint prediction of pixel and context using duplex neural networks. *IEEE Access*, 2021. 1
- [40] Oren Rippel and Lubomir Bourdev. Real-time adaptive image compression. In *International Conference on Machine Learning*, pages 2922–2930. PMLR, 2017. 1
- [41] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017. 1, 2
- [42] Ionut Schiopu and Adrian Munteanu. Deep-learning-based lossless image coding. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(7):1829–1842, 2019. 1
- [43] Jon Sneyers and Pieter Wuille. Flif: Free lossless image format based on maniac compression. In *2016 IEEE international conference on image processing (ICIP)*, pages 66–70. IEEE, 2016. 6
- [44] Myungseo Song, Jinyoung Choi, and Bohyung Han. Variable-rate deep image compression through spatially-adaptive feature transform. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2380–2389, 2021. 1
- [45] Yannick Strümpfer, Ren Yang, and Radu Timofte. Learning to improve image compression without changing the standard decoder. In *European Conference on Computer Vision*, pages 200–216. Springer, 2020. 1
- [46] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017. 1
- [47] George Toderici, Sean M O’Malley, Sung Jin Hwang, Damien Vincent, David Minnen, Shumeet Baluja, Michele Covell, and Rahul Sukthankar. Variable rate image compression with recurrent neural networks. *arXiv preprint arXiv:1511.06085*, 2015. 1
- [48] George Toderici, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor, and Michele Covell. Full resolution image compression with recurrent neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5306–5314, 2017. 1
- [49] Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International Conference on Machine Learning*, pages 1747–1756. PMLR, 2016. 1, 2
- [50] WebEngines Blazer Platform Version. 1.0 hardware reference guide, xp-002202892, network engines. Inc., Jun, 1:92, 2000. 6
- [51] Marcelo J Weinberger, Gadiel Seroussi, and Guillermo Sapiro. The loco-i lossless image compression algorithm: Principles and standardization into jpeg-ls. *IEEE Transactions on Image processing*, 9(8):1309–1324, 2000. 1, 6
- [52] Workshop and challenge on learned image compression. <https://www.compression.cc/challenge/>. 5
- [53] Fei Yang, Luis Herranz, Yongmei Cheng, and Mikhail G Mozerov. Slimmable compressive autoencoders for practical neural image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4998–5007, 2021. 1
- [54] Xi Zhang and Xiaolin Wu. Attention-guided image compression by deep reconstruction of compressive sensed saliency skeleton. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13354–13364, 2021. 1