

# High-Fidelity Variable-Rate Image Compression via Invertible Activation Transformation

Shilv Cai  
Huazhong University of Science and  
Technology  
caishilv@hust.edu.cn

Zhijun Zhang  
Huazhong University of Science and  
Technology  
zhangzhijun@hust.edu.cn

Liqun Chen  
Huazhong University of Science and  
Technology  
chenliqun@hust.edu.cn

Luxin Yan  
Huazhong University of Science and  
Technology  
yanluxin@hust.edu.cn

Sheng Zhong  
Huazhong University of Science and  
Technology  
zhongsheng@hust.edu.cn

Xu Zou  
Huazhong University of Science and  
Technology  
zoux@hust.edu.cn

## ABSTRACT

Learning-based methods have effectively promoted the community of image compression. Meanwhile, variational autoencoder (VAE) based variable-rate approaches have recently gained much attention to avoid the usage of a set of different networks for various compression rates. Despite the remarkable performance that has been achieved, these approaches would be readily corrupted once multiple compression/decompression operations are executed, resulting in the fact that image quality would be tremendously dropped and strong artifacts would appear (see Figure 1). Thus, we try to tackle the issue of high-fidelity fine variable-rate image compression and propose the Invertible Activation Transformation (IAT) module. We implement the IAT in a mathematical invertible manner on a single rate Invertible Neural Network (INN) based model and the quality level (QLevel) would be fed into the IAT to generate scaling and bias tensors. IAT and QLevel together give the image compression model the ability of fine variable-rate control while better maintaining the image fidelity. Extensive experiments demonstrate that the single rate image compression model equipped with our IAT module has the ability to achieve variable-rate control without any compromise. And our IAT-embedded model obtains comparable rate-distortion performance with recent learning-based image compression methods. Furthermore, our method outperforms the state-of-the-art variable-rate image compression method by a large margin, especially after multiple re-encodings.

## KEYWORDS

Image Compression; Variable-Rate; Fidelity Maintenance

### ACM Reference Format:

Shilv Cai, Zhijun Zhang, Liqun Chen, Luxin Yan, Sheng Zhong, and Xu Zou. 2022. High-Fidelity Variable-Rate Image Compression via Invertible Activation Transformation. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '22)*, October 10–14, 2022, Lisboa, Portugal. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3503161.3547880>

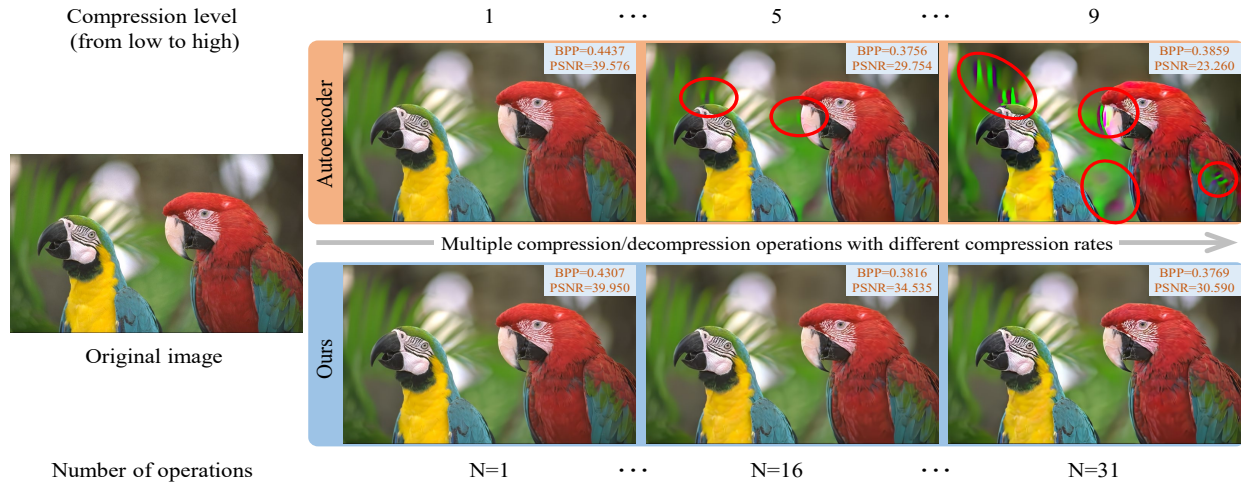
## 1 INTRODUCTION

Lossy image compression is one crucial technology due to the increasing volume of visual data in such a multimedia explosion era. This task aims at lowering data redundancy while maintaining visual fidelity and supporting efficient data storage and transmissions. To this end, many classical image compression standards (e.g., JPEG [49], JPEG2000 [41], Webp [17], BPG [7], and Versatile Video

Coding (VVC) [26]) have been proposed and widely used in practical applications. Recently, learning-based image compression methods have started to show superiority in terms of common metrics, e.g. PSNR and MS-SSIM. These methods [4, 5, 45] make use of the powerful nonlinear transformation capability of DNNs, and perform end-to-end learning by a large number of high-quality images with a rate-distortion cost. However, despite the exciting progress, the learning-based image compression still remains challenging once variable-rate compression adaptation is needed. Most of them require training multiple single-rate models for different rates, resulting in a high cost of model storage and training.

To remedy the issue, a large number of VAE-based variable-rate image compression methods [9, 11, 13, 43, 44, 55] have been proposed. The researchers first try to achieve discrete rate adaptation using one single model. Choi et al. [11] introduced conditional convolution and achieved variable rate through two-stage training. Yang et al. [55] proposed the modulated autoencoder and achieved discrete adjustable compression rates by different Lagrange multipliers. Chen et al. [9] inserted a set of scaling factors directly before the quantizer to achieve the discrete adjustable compression rates. However, the performance of these methods would be dropped when conducting finer variable-rate control. Thus, the topic of fine rate adaptation has attracted more attention recently. Sun et al. [44] obtained continuously adjustable compression rate by linear interpolation. Cui et al. [13] achieved continuous compression rate control by exponential interpolation. Song et al. [43] conditioned on quality map and achieved the variable rate, which requires semantic segmentation labels for training. Though these methods have the ability of fine variable-rate compression control, they need additional gain modules or semantic labels to maintain the performance.

Besides, it would be particularly interesting for a variable-rate compression model if the fidelity of images could be maintained while being transmitted multiple times between numerous entities under various compression rates, especially in the current multimedia society (e.g. one person may download a compressed image from Instagram and then send it to his friend via WhatsApp under another re-encoding). However, state-of-the-art VAE-based variable-rate approaches (e.g. Song et al. [43]) would be readily corrupted once multiple compression/decompression operations are executed, resulting in the fact that image quality would be tremendously dropped. Strong artifacts and color shifts would appear, as shown



**Figure 1: Reconstructed images of variable-rate compression methods after different numbers of compression/decompression operations. It broadly occurs in multimedia transmissions among social platforms. Severe artifacts and color shifts would appear (see regions in red circles) in the state-of-the-art VAE-based approach [43] once multiple re-encodings are executed, in contrast to fewer artifacts and higher fidelity results achieved by our proposed approach. High-fidelity maintenance with fine variable-rate control is the main advantage and novelty of our work.**

in Figure 1. The main reason is that the autoencoder transforms the image to a low-dimensional latent space and irreversibly discards information before quantization, imposing an implicit limitation on the reconstruction quality. To alleviate information loss, Invertible Neural Networks [19, 54] have gained much attention to effectively preserve fidelity. It is worth noting that VAE-based variable-rate methods cannot be directly fused into the INN-based framework since implementations of their variable-rate control do not satisfy the bijective mapping property. Nevertheless, there is no research on INN-based variable-rate methods to the best of our knowledge. Inspired by this, we construct a variable-rate image compression model which can maintain the fidelity, especially after multiple re-encodings, by exploring the invertibility.

To sum up, we propose an Invertible Activation Transformation (IAT) module based on the INN framework. This module exhibits a mathematical invertible property to avoid discarding any information in the latent space to maintain high fidelity. Notice that it is the initial work to extend the mathematical invertibility to the variable-rate image compression. Moreover, the proposed image compression method attempts to achieve finer control of multiple variable rates, by presenting a compatible tensor-based Lagrange multiplier to train the whole model. The contributions of our proposed method are 3-folded:

- We propose an effective yet neat framework, equipped with the INN-based Invertible Activation Transformation (IAT) module, to achieve the high fidelity of reconstructed images, especially after multiple variable-rate image compression/decompression operations, in a mathematical invertible manner. This issue is rarely investigated so far.
- The proposed model tuned rate-distortion loss and achieved fine variable-rate control through the quality level.
- Extensive experiments demonstrate the superiority of our proposed methods in rate-distortion performance, fidelity

maintenance, and fine rate adaptation over three datasets, including Kodak [12], CLIC [47], and DIV2K [1].

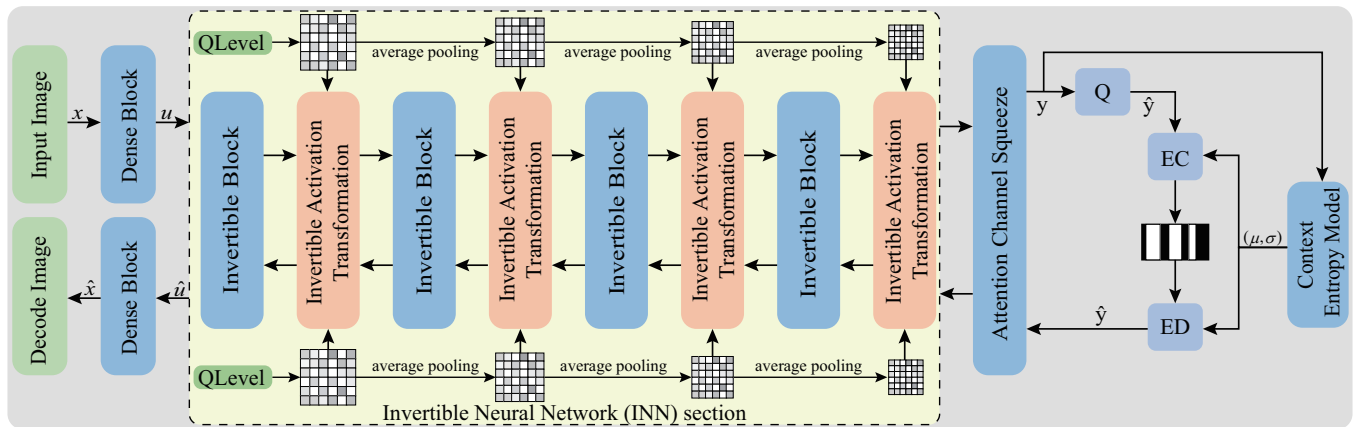
## 2 RELATED WORK

In recent years, the application of neural networks in image compression has attracted widespread attention. The variational autoencoder (VAE) [4, 5, 8, 10, 18, 21, 22, 29, 34, 35, 37, 38, 45, 57], Invertible Neural Network (INN) [19, 20, 33, 54] and Generative Adversarial Networks (GAN) [2, 24, 36, 42, 52] based methods have achieved surprising results.

### 2.1 Learned Single Rate Image Compression

The VAE-based framework is used as a nonlinear transformation coding model, which is the main approach in the learned image compression method. The works [4, 5, 45] were the first to use CNN for end-to-end image compression and inspired many learning-based image compression methods. The work [5] introduced a hyperprior entropy model to capture the zero-mean Gaussian distribution of the latent representations. The works [29, 37] used the Gaussian model with the non-zero mean to improve the ability to model latent representations. Later works [29, 35, 37] further removed redundancy in potential features using the context model. Further, the 3D-context entropy model [18], multi-scale hyperprior entropy model [22], and discretized Gaussian mixture model [10] were used to further improve the entropy model. In addition, channel-wise module [38], attention module [10, 57], and non-local attention module [8, 56] were used to extract better latent representations. Recently transformer was used to capture long-range dependencies in probability distribution estimation effectively and efficiently [40].

Most learning-based image compression methods need to train different network models for various compression rates, which not only increases the storage computational resources but also is not



**Figure 2: Network architecture equipped with the proposed Invertible Activation Transformation (IAT) module. We insert IAT into the Invertible Neural Network section and utilize it to generate element-wise activation parameters of features from the input quality level (QLevel). IAT and QLevel together give the image compression model the ability of fine variable-rate control while maintaining the image fidelity especially when multiple compression/decompression operations are executed. EC/ED means entropy encoding/decoding respectively. Q is the quantizer.**

compatible with practical applications. Therefore, using one single model to achieve variable rate adaptation was widely studied.

## 2.2 Learned Variable Rate Image Compression

Initially, LSTM networks [25, 46, 48] control different compression rates by the different number of iterations. The more iterations, the clearer the reconstructed image would be. However, the LSTM-based approach cannot outperform JPEG2000 [41] in rate-distortion performance and would not obtain continuous compression rates. In addition, the iterative procedure is very time-consuming and thus not suitable for practical applications. Then Choi et al. [11] introduced conditional convolution in the autoencoder framework to achieve variable-rate adaptation with a single model through two-stage training.

However, while variable rate is achieved, the rate-distortion performance degrades and there is a dilemma in choosing the appropriate Lagrange multiplier and quantization step size for forward inference. Yang et al. [55] proposed a modulated autoencoder that achieves discrete adjustable compression rate with a single model by different Lagrange multipliers. Thesis et al. [45] first trained the model with high bits per pixel (bpp) and then fixed the network model parameters to train the scaling parameters for different compression rates. However, the network model suffered from incongruity with the scaling parameters, especially in low bpp cases. Chen et al. [9] inserted a set of scaling factors directly before the quantizer to achieve the discrete variable compression rate.

Recently, research has been conducted on continuous compression rate adjustable [13, 43, 44]. The work [13] introduced a series of vector pairs for coarse compression rate control, and then achieve continuous compression rate control by exponential interpolation. Sun et al. [44] extended the work [11], which obtained a continuously adjustable compression rate by linear interpolation. Song et al. [43] conditioned the quality map by spatial feature transform (SFT) [51] to control different compression rates.

VAE-based variable-rate approaches have been extensively researched. However, those methods suffer from severe information distortion after multiple operations of compression/decompression for the same image. The distortion becomes more explicit as the number of operations increases.

## 2.3 Invertible Neural Networks

Invertible neural networks (INNs) are generative models that transform complex distributions into simple ones, allowing for accurate and efficient probability density estimation. INNs have a bijective mapping of input and output, which is ideal for image compression.

NICE [14] introduced a flexible architecture that can learn highly nonlinear bijective transformations to represent data with simple distributions. Based on NICE [14], RealNVP [15] further extended the idea of hierarchical and combinatorial transformations, which used affine coupling and a multi-scale framework. Kingma et al. [28] proposed a generative flow model based on a  $1 \times 1$  invertible convolutional network with a significant improvement in log-likelihood on a standard benchmark dataset, having the advantages of exact controllability of log-likelihood, the tractability of exact inference of latent representations, and parallelizability of training and synthesis. Ardizzone et al. [3] demonstrated that the validity of INNs is suitable not only for synthetic data but also for two practical applications in medicine and astrophysics. SRFlow [32] has designed a conditional normalizing flow architecture to solve the ill-posed problem in the super-resolution task. Xiao et al. [53] proposed an invertible rescaling network (IRN), which constructs a bijective transform to effectively implement the reconstruction of low-resolution images into high-resolution images.

INN greatly alleviates the information loss problem for better image compression, as in [19, 20, 33, 54]. But no one has specifically studied variable-rate image compression with a single model based on the INN framework.

### 3 METHODOLOGY

#### 3.1 Framework

Our image compression approach is depicted in Figure 2. The proposed method implements fine variable-rate modulation in an invertible neural network framework, which involves the invertible activation transformation (IAT) module to control different compression rates through different quality levels. We present the detailed procedure of the model in the following: Firstly, the source image  $x \in \mathbb{R}^{3 \times H \times W}$  is enhanced by the dense block module [23] to generate a nonlinear representation of  $u \in \mathbb{R}^{3 \times H \times W}$ , where  $H$  and  $W$  denote the height and width of the input image respectively. Then the forward pass of the Invertible Neural Network section, which is equipped with the proposed IAT module, transforms  $u$  to a latent representation, conditioned on the quality level  $L \in \mathbb{R}^{H \times W}$  to control the compression rate. This latent representation would be further fed into the Attention Channel Squeeze module to reduce the number of channels and obtain the potential representation  $y$ . This procedure could be formulated by a parametric analysis transform function, *i.e.*,

$$y = g_a(x, L), \quad (1)$$

the discrete latent features  $\hat{y}$  are obtained by quantification of  $y$ , *i.e.*,  $\hat{y} = Q(y)$ . We use the quantizer  $Q(\cdot)$  in Ballé et al. [5] to model the quantized latent representation  $\hat{y}$  approximately by adding the uniform noise  $U(-0.5, 0.5)$  to the latent representation  $y$  during training and rounding the latent representation  $y$  during testing. The context entropy model generates parameters  $\mu$  and  $\sigma$  of the Gaussian entropy model that approximates the distribution of quantified latent representation  $\hat{y}$  to support the entropy encoding. We use range asymmetric numeral system [16] to losslessly compress latent representation  $\hat{y}$  and  $\hat{z}$  into bitstreams.

The inverse calculation takes the quantified latent representation  $\hat{y}$  and the quality level  $L$  as the input, and reconstructs the decompressed images by a parametric synthesis transform, which is formulated as follows:

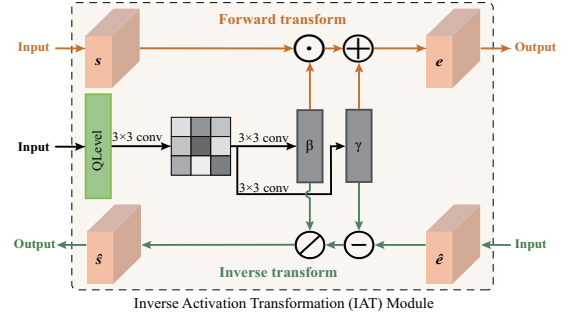
$$\hat{x} = g_s(\hat{y}, L). \quad (2)$$

#### 3.2 Invertible Activation Transformation

We proposed the invertible activation transformation (IAT) module to enhance the invertible neural network, which efficiently generates the desired compressed representation conditional on the quality level  $L$ . The proposed IAT module can achieve variable-rate adaption on a single model while maintaining the image fidelity, especially after multiple compression/decompression operations, in a mathematical invertible manner.

The forward transform of the IAT module is illustrated by pink arrows on the top of Figure 3. The inputs are the quality level  $L$  and the feature  $s$ . The element-wise activation parameters  $\gamma \in \mathbb{R}^{c \times h \times w}$  and  $\beta \in \mathbb{R}^{c \times h \times w}$  are then calculated by the IAT module from the quality level  $L$  via convolutional operations. These activation parameters would be applied to the feature  $s$  via the Equation 3 to generate the feature  $e$ ,

$$e = (s \odot \beta) \oplus \gamma, \quad (3)$$



**Figure 3: Illustration of the IAT module. The forward and inverse transformation of the IAT module implements the bijective mapping. This module takes the QLevel and feature as input to generate element-wise activation parameters  $\beta$  and  $\gamma$ , further obtaining the output results. Thus, the forward and inverse procedures are mathematically invertible, enhancing the fidelity of reconstructed images.**

where  $\odot$  denotes the Hadamard product and  $\oplus$  denotes the addition by element.  $c$ ,  $h$ , and  $w$  are the channel, height, and width of the feature, respectively.

The inverse transform of the IAT module is illustrated by green arrows at the bottom of Figure 3. The input quality level  $L$  and features  $\hat{e}$  are applied to obtain the feature  $\hat{s}$ . This inverse transform is formulated by Equation 4,

$$\hat{s} = (\hat{e} \odot \gamma) \oslash \beta, \quad (4)$$

where  $\ominus$  denotes the subtraction in elemental order,  $\oslash$  denotes the division by elemental order. Once the quality level  $L$  is the same in both forward and inverse procedures, the invertibility of the operation between the features  $s$  and  $e$  can be guaranteed.

In the previous work [9], a set of scaling factors was inserted directly before the quantizer to achieve the discrete adjustable compression rate. In our algorithm, the activation parameters are element-wise, which means that IAT module is computed as a spatial feature transform rather than a simple channel weighting. Moreover, the IAT module is attached after each invertible block which is initially proposed in RealNVP [15] and adopted by baseline model [54], not just inserted before the quantizer. These adjustments not only make fine variable-rate adaptation available but also turn out to the better performance, the experiment "Impact of the QLevel Representation" in section 4.4 shows its effectiveness, and the results are shown in Figure 7.

#### 3.3 Fine Variable-Rate Control

Unlike interpolation-based methods [13, 44] for obtaining finer compression rates, our method achieves the fine compression rate adaptation directly by modulating the quality level  $L$ , which is more convenient when controlling the compression rate by only one parameter instead of two. Compared to Song et al. [43], our method does not require additional semantic labels, either.

The goal of lossy image compression is to minimize the length of the bits stream and the distortion between the source image  $x$  and the reconstructed image  $\hat{x}$ . The optimization function is always expressed in the rate-distortion loss:  $R + \lambda D$ , where  $\lambda$  is

the Lagrange multiplier which determines the trade-off between the rate  $R$  and the distortion  $D$ . Theoretically, as long as the set of Lagrangian multiplier  $\lambda$  is large enough, it is possible to achieve fine compression rate control, but in practice, the computational cost is too high. For interpolation-based methods, the Lagrangian multiplier  $\lambda$  is a scalar. Thus, at each iteration during training, only one element in a finite set of  $\lambda$  would be randomly selected for optimization. In order to further promote the R-D performance of our model, we use a tensor instead of the scalar  $\lambda$ . Our optimization function implements fine variable-rate control by minimizing the rate-distortion loss  $R + \Lambda \odot \mathbf{D}$ , where dimensions of  $\Lambda \in \mathbb{R}^{C \times H \times W}$  and the distortion  $\mathbf{D} \in \mathbb{R}^{C \times H \times W}$  are the same as the dimension of the original input image.  $\odot$  denotes the Hadamard product. In this formulation,  $\Lambda$  is a tensor and no longer a finite set of constant scalars. Thus,  $\mathbf{D}$  measures pixel-wise distortion and is defined as  $\mathbf{D} = \frac{\sum_{i=1}^T \lambda_i (x_i - \hat{x}_i)^2}{T}$ ,  $T$  indicates the number of image pixels,  $\lambda_i$  is the Lagrangian multiplier,  $x_i$  and  $\hat{x}_i$  denote one pixel of the original and reconstructed image, respectively.

$\Lambda$  is simply calculated from the quality level  $L$  via a monotonically increasing function:  $\Lambda = V(L)$ , where  $V : \mathbb{R}^N \rightarrow \mathbb{R}^T$ .  $V(L) = \theta \times e^{\tau \times L}$ ,  $\theta = 0.0012$ ,  $\tau = 4.382$ , the process of dimensioning from  $\mathbb{R}^N \rightarrow \mathbb{R}^T$  is done by direct replication between channels.  $L = [l_i]_{i=1:N}$ ,  $l_i \in [0, 1]$ ,  $N = H \times W$ ,  $T = C \times H \times W$ .  $C$ ,  $H$ , and  $W$  denote the channel, height, and width of the source image  $x$ , respectively. Under such a paradigm, we implement this pixel-wise distortion constraint by randomly generating values of each element of the tensor  $\Lambda$  via the quality level  $L$  during training. This is equivalent to increasing the number of  $\lambda$  values selected at each iteration. So, the fine variable-rate control can be obtained by feeding exact quality levels during the testing.

As in other learning-based method [5], the log-likelihood of the coded features  $\hat{y}$  is estimated by a probabilistic model to replace the true compression rate  $R$ . Finally, the training loss would be:

$$\text{Loss} = -\log_2 P_{\hat{y}}(\hat{y}|\Lambda) - \log_2 P_{\hat{z}}(\hat{z}|\Lambda) + \frac{\sum_{i=1}^T \lambda_i (x_i - \hat{x}_i)^2}{T}, \quad (5)$$

where  $\hat{y}$  and  $\hat{z}$  are quantized latent representations and side information respectively.  $p_{\hat{y}}(\hat{y}|\Lambda) = \mathcal{N}(\mu, \sigma^2)$ ,  $\mu$  and  $\sigma$  denote the estimates of the mean and standard deviation of the quantified latent representation  $\hat{y}$ .  $p_{\hat{z}}(\hat{z}|\Lambda) = \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $\mu_1$  and  $\sigma_1$  denote the estimates of the mean and standard deviation of the quantified side information  $\hat{z}$ . The side information usually represents the hyper-prior originally proposed in [5] and refers to the extra stream  $\hat{z}$  generated by the "Context Entropy Model" in Figure 2. It is worth noting that this loss function would be degraded to the standard rate-distortion optimization function if all elements of the tensor quality level  $L$  are the same.

In addition, our method can be trained on arbitrary unlabeled data instead of requiring semantic segmentation labels corresponding to the original data, which is different from Song et al. [43], for training the model.

## 4 EXPERIMENTS

### 4.1 Implementation Details

**Details For Training.** In our implementation, the network of Xie et al. [54] is adopted as our basic architecture. The training datasets

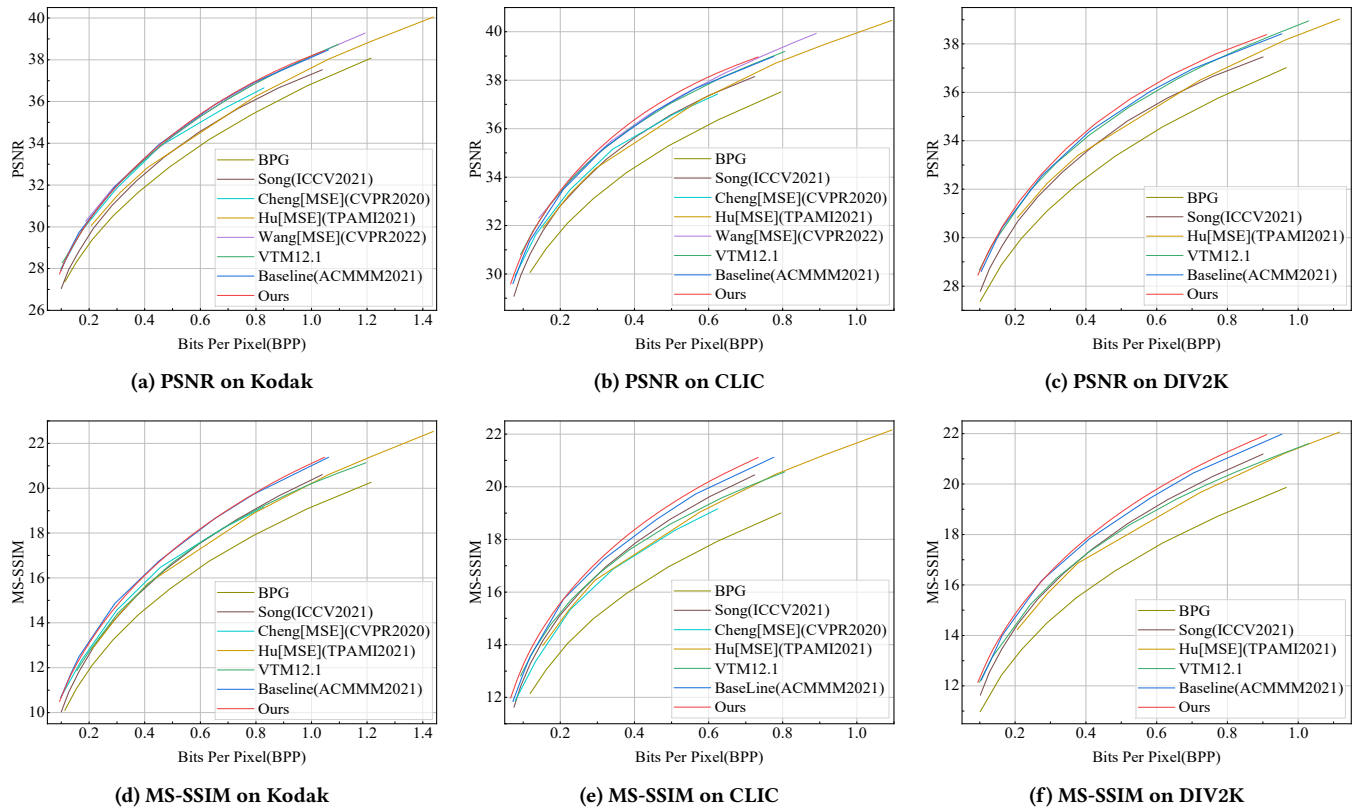
contain Flicker 2W [31] and COCO [30]. Flicker is used to train the network which has the context model, COCO is used to train the network without the context model. Our network is trained on  $256 \times 256$  randomly cropped patches and discards images less than 256px in height or width during data pre-processing. In training, the quality level  $L$  needs to be sent to the INN section as a condition during the forward and inverse transform. The quality level  $L$  takes a uniform value tensor between (0,1) during the testing and is randomly sampled between (0,1) during the training. Our implementation relies on Pytorch [39] and an open-source CompressAI PyTorch library [6]. All experiments were conducted on RTX 3090 GPU and trained for about 2.5M iterations with batch size 8. Adam optimizer [27] is used to optimize the parameters, there were multistage learning rates  $\{1e-4, 5e-5, 1e-5, 5e-6, 1e-6, 5e-7\}$  that changed with boundaries  $\{1000000, 1300000, 1600000, 1900000, 2200000, 2500000\}$ .

**Details For Testing.** We evaluate the rate-distortion performance on three commonly used datasets. The Kodak [12] contains 24 lossless images with a size of  $768 \times 512$ . The CLIC Professional Validation dataset [47] comprises 41 high-quality images with much higher resolution. The DIV2K validation dataset [1] contains 100 images with high resolutions of 2K. We draw curves based on the rate-distortion performance to compare the coding efficiency of different methods. We also calculate the area under the rate-distortion curve to observe the performance difference more effectively.

### 4.2 Rate-Distortion Performance

To verify the validity of the proposed approach, we conduct rate-distortion (RD) performance experiments on three datasets, *i.e.*, Kodak [12], CLIC [47], and DIV2K [1]. We compare our approach with five recent state-of-the-art learning-based image compression methods [10, 22, 43, 50, 54] and two classical codec methods, BPG [7] and VCC [26]. The results of learning-based methods are collected from their official GitHub pages or their papers. The VCC approach is implemented by the official Test Model VTM 12.1 with the intra-profile configuration from the official GitHub page to test images. Both VVC and BPG software were configured with the YUV444 format to maximize compression performance.

All comparable results are demonstrated in Figure 4. It is seen that our approach achieves the best results with commonly used metrics PSNR and MS-SSIM on three datasets. Compared with the baseline method [54], our approach achieves comparable R-D performance on the Kodak dataset (Figure 4 (a)(d)) and outperforms the baseline on both the CLIC dataset (Figure 4 (b)(e)) and the DIV2K dataset (Figure 4 (c)(f)). This means that our approach achieves the variable-rate adaptation based on the single rate method [54] without sacrificing any performance, verifying the effectiveness of the IAT module. It is worth noting that the CLIC dataset and DIV2K dataset are high-resolution images, implying that our method is more effective on high-resolution images. Our approach empowers the network model with variable-rate in addition to improving the algorithmic performance of the original model. To further compare the performance between the baseline [54] and our method, we calculate their corresponding area under curve (AUC) values, as shown in Table 1. The results show that our approach outperforms the single rate model method by Xie et al. [54] in terms of the aggregated AUC metric.



**Figure 4: RD performance curves aggregated over the Kodak, CLIC professional validation dataset, and DIV2K validation dataset. MS-SSIM values converted to decibels ( $-10\log_{10}(1 - MS - SSIM)$ ). (a)/(b)/(c) and (d)/(e)/(f) are results on Kodak, CLIC, and DIV2K about PSNR and MS-SSIM, respectively. It is worth noting that CLIC and DIV2K are datasets with high-resolution images. That is, our method is especially effective on high-resolution images.**

In addition, Our proposed method could achieve variable-rate image compression with a fine granularity. To verify the effectiveness of fine variable-rate control, we illustrate multiple performances of fine variable-rate control within the low and high bpp range in Table 2. In practice, classical image codecs provide hundreds of variable-rate RD points to meet the basic requirement of the application. Compared with that, our method obtains at least 1000 effective variable-rate RD points with a very fine PSNR and MS-SSIM. We achieved the fine-rate control compared with the classical image codecs BPG [7] and VTM 12.1 [26], the comparative results refer to the supplementary material.

### 4.3 Fidelity for Re-encoding

In order to verify the high fidelity, our method is compared with the latest VAE-based variable-rate method by Song et al. [43]. This method does not use context model and has available source code. To make a fair comparison, we remove the context model and add the non-local attention module [8] to the hyperprior layer. Figure 5 illustrates the results of multiple compression/decompression operations on the same image with different compression rates. With operations increasing, our proposed method shows higher fidelity.

**Table 1: Area under curve (AUC) of our method and Xie et al. [54](Baseline) on different datasets about PSNR and MS-SSIM. The bpp range is determined by the intersection of two methods. Our approach makes a single-rate baseline compression model achieve the variable-rate ability and even outperforms the baseline in R-D performance.**

Dataset	Xie et al. [54]		Ours	
	$AUC_{PSNR}$	$AUC_{MS-SSIM}$	$AUC_{PSNR}$	$AUC_{MS-SSIM}$
Kodak	32.7866	16.5030	<b>32.7883</b>	<b>16.5036</b>
CLIC	23.5896	11.7463	<b>23.7082</b>	<b>11.8571</b>
DIV2K	28.0998	14.7868	<b>28.2138</b>	<b>14.8901</b>

Figure 6 (a)(c) show the rate-distortion performance after multiple operations of compression/decompression with different compression rates. Both approaches change from high to low bpp ranges, our method in the set of bpp { 1.0267, 1.0116, 0.9949, 0.9784, 0.9619, 0.9456, 0.9292, 0.9127, 0.8965, 0.8809, 0.8658, 0.8507, 0.8357, 0.8206, 0.8056, 0.7907 }, Song et al. [43] in the set of bpp {1.0392, 1.0249, 1.0091, 0.9932, 0.9768, 0.9606, 0.9449, 0.9287, 0.9128, 0.8968, 0.8813, 0.8658, 0.8505, 0.8351, 0.8201, 0.8052}. It is clearly seen

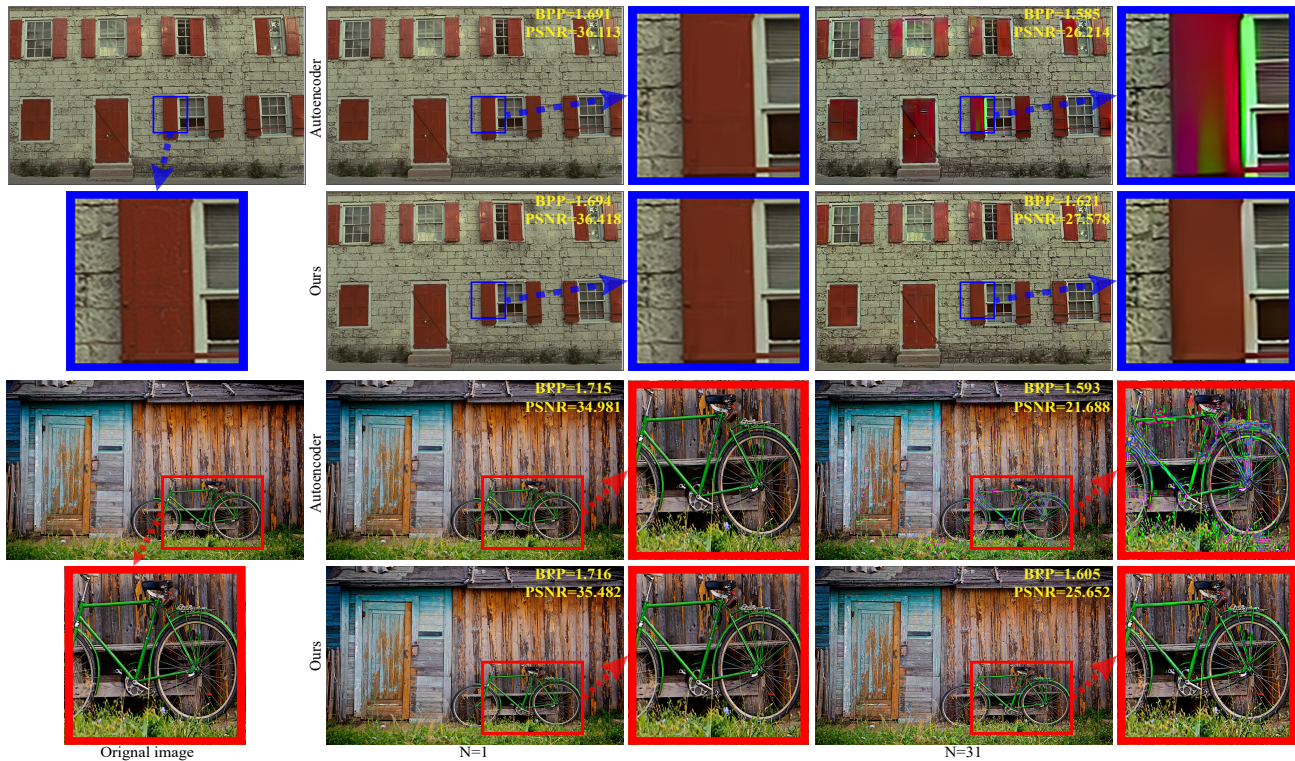


Figure 5: Qualitative results after different numbers of compression/decompression operations under various rates. The two images (kodim1.png and alexander-shustov-73.png) are from the Kodak dataset and the CLIC dataset, respectively. Severe artifacts and color shifts would appear in the state-of-the-art VAE-based approach [43] once multiple operations are executed, in contrast to better fidelity maintenance of our approach. Please refer to the supplementary material for more cases. N indicates the number of compression/decompression operations. Best viewed in color.

Table 2: Variable-rate control experiments over the Kodak dataset. Our approach can finely control the compression rate within the whole bpp range (no matter low or high).

LOW			HIGH		
BPP	PSNR(dB)	MS-SSIM(dB)	BPP	PSNR(dB)	MS-SSIM(dB)
0.28181	31.6951	14.5015	1.02433	38.3226	21.2580
0.28265	31.7071	14.5153	1.02587	38.3312	21.2664
0.28342	31.7177	14.5263	1.02733	38.3388	21.2717
0.28416	31.7291	14.5377	1.02910	38.3468	21.2819
0.28500	31.7435	14.5517	1.03071	38.3548	21.2903
0.28576	31.7538	14.5639	1.03250	38.3625	21.2995
0.28659	31.7657	14.5765	1.03406	38.3703	21.3087
0.28734	31.7761	14.5874	1.03564	38.3767	21.3190
0.28808	31.7884	14.5952	1.03733	38.3872	21.3291
0.28880	31.8004	14.6092	1.03885	38.3943	21.3355

that our method outperforms Song et al. [43], after multiple compression/decompression operations. Figure 6 (b)(d) show the rate-distortion performance by multiple operations with the fixed compression rate. Both approaches achieve a bit rate of 0.791 bpp for all steps. Also, our method achieves better results significantly, compared with Song et al. [43] and baseline [54]. The results indicate

that our proposed IAT module is powerful to maintain image fidelity, which is important for practical applications. It is noteworthy that compression methods capable of high-fidelity in re-encoding are of great importance in the video production pipeline, as image/video content may be edited/composited by different people or at different times, requiring re-encodings in the process.

#### 4.4 Discussion

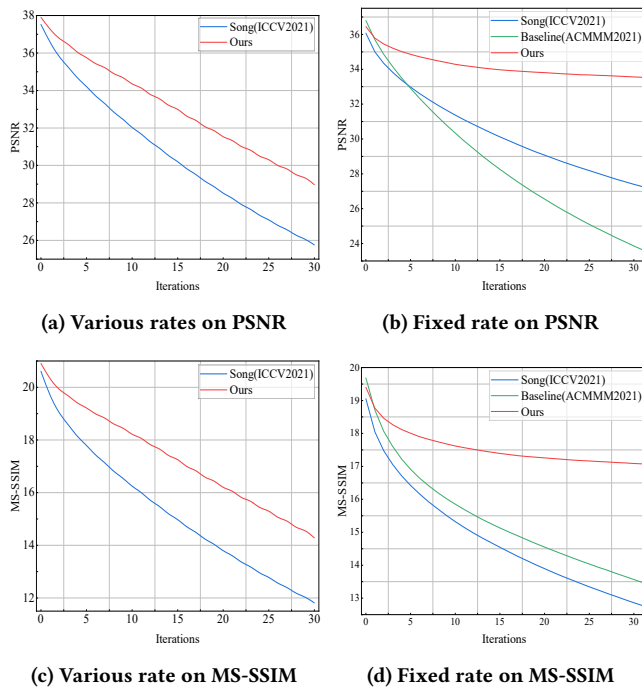
**Impact of the QLevel Representation.** To further analyze the effectiveness of the tensor-based QLevel representation of our IAT module, we conducted an ablation study by modifying the quality level representation. We compared our approach with the baseline method [54] and the simplified version of our method, which modifies the quality level from tensor to scalar, similar to [9]. Comparative results are shown in Figure 7. The results indicate that the proposed tensor-based quality level obtains better performance, compared with the scalar factor ones, which only provides channel-wise weighted computations on latent representation.

**Impact of Gain Components.** The context model [29, 35, 37] and the non-local attention module [8] are commonly used in the learned-based image compression methods to further reduce statistical redundancy within the latent features and improve the probabilistic estimation ability of the network. We conduct an ablation

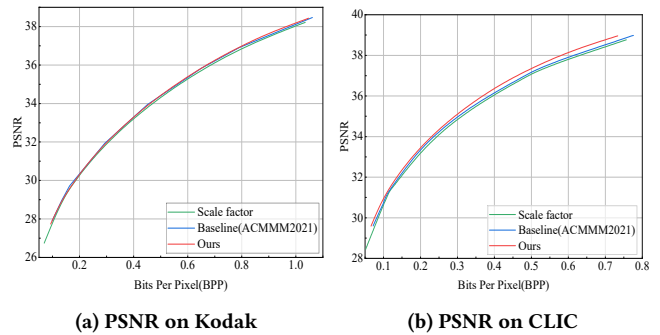
study to evaluate the impact of the context model and non-local attention module on our method in the Kodak dataset, as shown in Figure 8. We start from a baseline without the context model and non-local attention module, *i.e.*, W/O CM (context model) and W/O NLAM (non-local attention module), and plot the rate-distortion performance in green color. Then, we add the non-local attention module (blue color) and context model (red color) to evaluate the performance. We can observe that using the context model achieves the best results, while it requires high computational costs (codec process takes about 233 seconds on an Intel(R) Core(TM) i9-10900K CPU). In addition, Our method outperforms Song et al. [43] without the context model and non-local attention module, demonstrating the effectiveness of the proposed method.

## 5 CONCLUSION

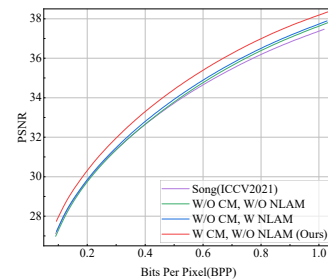
In this paper, we propose a high-fidelity variable-rate image compression method by introducing the Invertible Activation Transformation (IAT) module. The IAT module, implemented in a mathematical invertible manner, as a feature activation transform layer of the invertible neural network, has the ability of fine variable-rate control by feeding the quality level (QLevel) to generate the scaling



**Figure 6: Successive re-encodings on the Kodak dataset. (a) and (c): Compression rates of each compression/decompression operation are different. (b) and (d): The compression rate is fixed. Our approach outperforms baseline [54] and Song et al. [43] (a SOTA variable-rate approach) by a large margin to show the superiority of fidelity maintenance especially when multiple operations are executed.**



**Figure 7: Impact of the QLevel Representation. The scale factor method (green line) is similar to Chen et al. [9]. Our proposed tensor-based QLevel representation achieves better performance than simply using a scalar to control the compression rate.**



**Figure 8: Impact of Gain Components. W/O represents ‘without’, W represents ‘with’, CM represents ‘context model’, and NLAM represents ‘non-local attention module’.**

and bias tensors while better maintaining the image fidelity. Extensive experiments demonstrate that the single rate model equipped with our IAT module is able to achieve fine variable-rate control without any performance compromise. Thanks to the mathematical invertibility of our approach, fewer artifacts or color shifts would have appeared and the fidelity of reconstructed images is better maintained, especially when multiple re-encodings are executed under various compression rates.

## ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China (NSFC) grant No. 62176100, the Special Project of Science and Technology Development of Central guiding Local-Central Guidance on Local Science and Technology Development Fund of Hubei Province grant 2021BEE056 and the National Key Laboratory Foundation of China grant No. 6142113200307.



## REFERENCES

- [1] Eirikur Agustsson and Radu Timofte. 2017. NTIRE 2017 Challenge on Single Image Super-Resolution: Dataset and Study. In *CVPRW*.
- [2] Eirikur Agustsson, Michael Tschannen, Fabian Mentzer, Radu Timofte, and Luc Van Gool. 2019. Generative Adversarial Networks for Extreme Learned Image Compression. In *ICCV*.
- [3] Lynton Ardizzone, Jakob Kruse, Sebastian Wirkert, Daniel Rahner, Eric W Pellegrini, Ralf S Klessen, Lena Maier-Hein, Carsten Rother, and Ullrich Köthe. 2019. Analyzing Inverse Problems with Invertible Neural Networks. In *ICLR*.
- [4] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. 2017. End-to-End Optimized Image Compression. In *ICLR*.
- [5] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. 2018. Variational Image Compression with a Scale Hyperprior. In *ICLR*.
- [6] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja. 2020. CompressAI: a Pytorch Library and Evaluation Platform for End-to-End Compression Research. *arXiv preprint arXiv:2011.03029* (2020).
- [7] Fabrice Bellard. 2015. BPG Image Format. <https://bellard.org/bpg/>
- [8] Tong Chen, Haojie Liu, Zhan Ma, Qiu Shen, Xun Cao, and Yao Wang. 2021. End-to-End Learned Image Compression via Non-Local Attention Optimization and Improved Context Modeling. *IEEE Transactions on Image Processing* 30 (2021), 3179–3191.
- [9] Tong Chen and Zhan Ma. 2020. Variable Bitrate Image Compression with Quality Scaling Factors. In *ICASSP*.
- [10] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. 2020. Learned Image Compression with Discretized Gaussian Mixture Likelihoods and Attention Modules. In *CVPR*.
- [11] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. 2019. Variable Rate Deep Image Compression with a Conditional Autoencoder. In *ICCV*.
- [12] Eastman Kodak Company. 1999. Kodak Lossless True Color Image Suite. <http://r0k.us/graphics/kodak/>
- [13] Ze Cui, Jing Wang, Shangyin Gao, Tiansheng Guo, Yihui Feng, and Bo Bai. 2021. Asymmetric Gained Deep Image Compression with Continuous Rate Adaptation. In *CVPR*.
- [14] Laurent Dinh, David Krueger, and Yoshua Bengio. 2015. NICE: Non-Linear Independent Components Estimation. In *ICLRW*.
- [15] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. 2017. Density Estimation Using Real NVP. In *ICLR*.
- [16] Jarek Duda. 2013. Asymmetric Numeral Systems: Entropy Coding Combining Speed of Huffman Coding with Compression Rate of Arithmetic Coding. *arXiv preprint arXiv:1311.2540* (2013).
- [17] Google. 2010. Web Picture Format. <https://chromium.googlesource.com/webm/libwebp>
- [18] Zongyu Guo, Yaojun Wu, Runsen Feng, Zhizheng Zhang, and Zhibo Chen. 2020. 3-D Context Entropy Model for Improved Practical Image Compression. In *CVPRW*.
- [19] Leonhard Helminger, Abdelaziz Djelouah, Markus Gross, and Christopher Schroers. 2020. Lossy Image Compression with Normalizing Flows. In *CoRR*.
- [20] Yung-Han Ho, Chih-Chun Chan, Wen-Hsiao Peng, Hsueh-Ming Hang, and Marek Domański. 2021. ANFIC: Image Compression Using Augmented Normalizing Flows. *IEEE Open Journal of Circuits and Systems* 2 (2021), 613–626.
- [21] Yueyu Hu, Wenhan Yang, and Jiaying Liu. 2020. Coarse-to-Fine Hyper-Prior Modeling for Learned Image Compression. In *AAAI*.
- [22] Yueyu Hu, Wenhan Yang, Zhan Ma, and Jiaying Liu. 2021. Learning End-to-End Lossy Image Compression: A Benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [23] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely Connected Convolutional Networks. In *CVPR*.
- [24] Shoma Iwai, Tomo Miyazaki, Yoshihiro Sugaya, and Shinichiro Omachi. 2021. Fidelity-Controllable Extreme Image Compression with Generative Adversarial Networks. In *ICPR*.
- [25] Nick Johnston, Damien Vincent, David Minnen, Michele Covell, Saurabh Singh, Troy Chinen, Sung Jin Hwang, Joel Shor, and George Toderici. 2018. Improved Lossy Image Compression With Priming and Spatially Adaptive Bit Rates for Recurrent Networks. In *CVPR*.
- [26] Joint Video Experts Team (JVET). 2021. VVC Official Test Model VTM. [https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware\\_VTM/-/tree/VTM-12.1](https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM/-/tree/VTM-12.1) accessed on April 5, 2021.
- [27] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*.
- [28] Durk P Kingma and Prafulla Dhariwal. 2018. Glow: Generative Flow with Invertible 1x1 Convolutions. In *NeurIPS*.
- [29] Jooyoung Lee, Seunghyun Cho, and Seung-Kwon Beack. 2019. Context-Adaptive Entropy Model for End-to-End Optimized Image Compression. In *ICLR*.
- [30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft COCO: Common Objects in Context. In *ECCV*.
- [31] Jiaheng Liu, Guo Lu, Zhihao Hu, and Dong Xu. 2020. A Unified End-to-End Framework for Efficient Deep Image Compression. *arXiv preprint arXiv:2002.03370* (2020).
- [32] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. 2020. SRFflow: Learning the Super-Resolution Space with Normalizing Flow. In *ECCV*.
- [33] Haichuan Ma, Dong Liu, Ning Yan, Houqiang Li, and Feng Wu. 2020. End-to-End Optimized Versatile Image Compression with Wavelet-Like Transform. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44 (2020), 1247–1263.
- [34] Yi Ma, Yongqi Zhai, Jiayu Yang, Chunhui Yang, and Ronggang Wang. 2021. AFEC: Adaptive Feature Extraction Modules for Learned Image Compression. In *ACMMM*.
- [35] Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. 2018. Conditional Probability Models for Deep Image Compression. In *CVPR*.
- [36] Fabian Mentzer, George D Toderici, Michael Tschannen, and Eirikur Agustsson. 2020. High-Fidelity Generative Image Compression. In *NeurIPS*.
- [37] David Minnen, Johannes Ballé, and George D Toderici. 2018. Joint Autoregressive and Hierarchical Priors for Learned Image Compression. In *NeurIPS*.
- [38] David Minnen and Saurabh Singh. 2020. Channel-Wise Autoregressive Entropy Models for Learned Image Compression. In *ICIP*.
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*.
- [40] Yichen Qian, Ming Lin, Xiuyu Sun, Zhiyu Tan, and Rong Jin. 2022. Entroformer: A Transformer-Based Entropy Model for Learned Image Compression. In *ICLR*.
- [41] Majid Rabbani. 2002. JPEG2000: Image Compression Fundamentals, Standards and Practice. *Journal of Electronic Imaging* 11, 2 (2002), 286.
- [42] Oren Rippel and Lubomir Bourdev. 2017. Real-Time Adaptive Image Compression. In *ICML*.
- [43] Myungseo Song, Jinyoung Choi, and Bohyung Han. 2021. Variable-Rate Deep Image Compression through Spatially-Adaptive Feature Transform. In *ICCV*.
- [44] Zhenhong Sun, Zhiyu Tan, Xiuyu Sun, Fangyi Zhang, Yichen Qian, Dongyang Li, and Hao Li. 2021. Interpolation Variable Rate Image Compression. In *ACMMM*.
- [45] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. 2017. Lossy Image Compression with Compressive Autoencoders. In *ICLR*.
- [46] George Toderici, Sean M O'Malley, Sung Jin Hwang, Damien Vincent, David Minnen, Shumeet Baluja, Michele Covell, and Rahul Sukthankar. 2015. Variable Rate Image Compression with Recurrent Neural Networks. In *ICLR*.
- [47] George Toderici, Wenzhe Shi, Radu Timofte, Johannes Balle Lucas Theis, Eirikur Agustsson, Nick Johnston, and Fabian Mentzer. 2021. Workshop and Challenge on Learned Image Compression. <http://www.compression.cc>
- [48] George Toderici, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor, and Michele Covell. 2017. Full Resolution Image Compression with Recurrent Neural Networks. In *CVPR*.
- [49] Gregory K Wallace. 1992. The JPEG Still Picture Compression Standard. *IEEE Transactions On Consumer Electronics* 38, 1 (1992), 18–34.
- [50] Dezhao Wang, Wenhan Yang, Yueyu Hu, and Jiaying Liu. 2022. Neural Data-Dependent Transform for Learned Image Compression. *arXiv preprint arXiv:2203.04963* (2022).
- [51] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. 2018. Recovering Realistic Texture in Image Super-Resolution by Deep Spatial Feature Transform. In *CVPR*.
- [52] Lirong Wu, Kejie Huang, and Haibin Shen. 2020. A Gan-Based Tunable Image Compression System. In *WACV*.
- [53] Mingqing Xiao, Shuxin Zheng, Chang Liu, Yaolong Wang, Di He, Guolin Ke, Jiang Bian, Zhouchen Lin, and Tie-Yan Liu. 2020. Invertible Image Rescaling. In *ECCV*.
- [54] Yueqi Xie, Ka Leong Cheng, and Qifeng Chen. 2021. Enhanced Invertible Encoding for Learned Image Compression. In *ACMMM*.
- [55] Fei Yang, Luis Herranz, Joost Van De Weijer, José A Iglesias Guitián, Antonio M López, and Mikhail G Mozerov. 2020. Variable Rate Deep Image Compression with Modulated Autoencoder. *IEEE Signal Processing Letters* 27 (2020), 331–335.
- [56] Yulun Zhang, Kungpeng Li, Kai Li, Bineng Zhong, and Yun Fu. 2019. Residual Non-Local Attention Networks for Image Restoration. In *ICLR*.
- [57] Lei Zhou, Zhenhong Sun, Xiangji Wu, and Junmin Wu. 2019. End-to-End Optimized Image Compression with Attention Mechanism. In *CVPRW*.