# Towards Backward-Compatible Continual Learning of Image Compression

Zhihao Duan[1]    Ming Lu[2]    Justin Yang[1]    Jiangpeng He[1†]    Zhan Ma[2]    Fengqing Zhu[1]

[1] Purdue University, West Lafayette, Indiana, U.S.A.

[2] Nanjing University, Nanjing, Jiangsu, China

{duan90, yang1834, he416, zhu0}@purdue.edu, {minglu, mazhan}@nju.edu.cn

## Abstract

*This paper explores the possibility of extending the capability of pre-trained neural image compressors (e.g., adapting to new data or target bitrates) without breaking backward compatibility, the ability to decode bitstreams encoded by the original model. We refer to this problem as continual learning of image compression. Our initial findings show that baseline solutions, such as end-to-end fine-tuning, do not preserve the desired backward compatibility. To tackle this, we propose a knowledge replay training strategy that effectively addresses this issue. We also design a new model architecture that enables more effective continual learning than existing baselines. Experiments are conducted for two scenarios: data-incremental learning and rate-incremental learning. The main conclusion of this paper is that neural image compressors can be fine-tuned to achieve better performance (compared to their pre-trained version) on new data and rates without compromising backward compatibility. Our code is available at this link.*

## 1. Introduction

Recent years have witnessed the rapid development of deep learning-based image compression. Most existing research in this field considers the *offline learning* setting, *i.e.*, once a neural network model is trained, its parameters are fixed and kept unchanged when deployed. However, real-world applications are often complex and dynamic, and an ideal compressor should be capable of being incrementally learned and adapted to various scenarios. For example, consider an image storage application with a compressor pre-trained on natural images for certain predefined target bitrates. As new image sources (*e.g.*, microscopy, remote sensing, and human face images) are encountered, one may want to update the compressor to improve its performance on the new data and to support different target rates. This raises an interesting question: *can neural network-based image com-*



(a) Continual learning of image compression



Original image | Reconstruction, before fine-tuning | **Vanilla fine-tuning** | **With our method**

Reconstruction, after fine-tuning

(b) The backward compatibility problem. Once a neural compressor (in this experiment, [41]) is fine-tuned, it can no longer decode the bitstream produced by its original version. Our method addresses this issue.
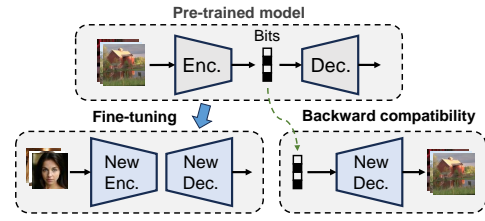
Figure 1. The goal of this work is to fine-tune pre-trained learned image compressors with new data or new rates while preserving backward compatibility (Fig. 1a). Baseline methods are backward incompatible, while ours is effective (Fig. 1b).

*pressors be learned continually, and if so, would it bring performance benefits compared to the pre-trained model?*

One might assume that simply fine-tuning a pre-trained model would be sufficient. Yet, doing so disrupts the model's *backward compatibility* (Fig. 1b), *i.e.*, the ability to decode bitstreams produced by the original model, due to a mismatch between the encoder (pre-trained model) and decoder (fine-tuned model). Maintaining backward compatibility is crucial, as failing to do so renders existing bitstreams in the storage (or sent from other devices) inaccessible. It is worth noting that this backward compatibility problem is different from the well-known problem of catastrophic forgetting in neural networks [28]. The unique properties of compression, including the sender-receiver relationship and the existence of entropy coding, set it apart from other image processing/vision tasks. Therefore, existing continual learning methods for vision tasks [13, 51] cannot be applied as is, and new strategies must be developed for compression to maintain the decoder's backward

---

[†]Corresponding author & Project lead

compatibility when adapting to new data and rates.

To achieve backward compatibility, the most straightforward way is to modify only the encoding process of a pre-trained model when adapting to new data or new rates [7, 19, 56]. This strategy resembles the common practice in traditional codecs, where there is often a standardized decoder, but various encoders can be designed to accommodate different applications. Despite its simplicity, keeping the decoder unchanged hinders the model's ability to adapt to new data and rates, leading to suboptimal performance.

This work shows that it is possible to continually train both the encoder and decoder of neural compressors while maintaining backward compatibility. We begin by noticing that as long as the entropy model of a compressor is kept unchanged, it can decode the old bitstreams and obtain the latent representations. Based on this observation, we propose a knowledge replay training scheme that can be used to train the encoder and decoder networks without breaking backward compatibility. We also design a model architecture where the entropy model consumes only a small amount of parameters, allowing most model parameters to be learnable in fine-tuning. We formulate two experimental scenarios: data-incremental learning and rate-incremental learning. Experimental results demonstrate that our proposed methods enable neural image compressors to obtain improved performance on new data and new rates without breaking backward compatibility.

To summarize, our contributions are as follows:

- We propose a knowledge replay-based training strategy that can be used to train neural image compressors incrementally without breaking backward compatibility;
- We design a neural network architecture targeting effective continual learning of image compression;
- We formulate two continual learning scenarios for image compression: data-incremental learning and rate-incremental learning. Experimental results show that our method outperforms baseline approaches in both cases.

## 2. Background and Related Work

### 2.1. Learned lossy image compression (LIC)

Most learning-based methods for lossy image compression can be interpreted using the *entropy-constrained non-linear transform coding* framework [4]. Let $X \sim p_{\text{data}}$ denote data samples with an underlying data distribution. In this framework, a neural network encoder $f_{\text{enc}}$ maps $X$ to a latent variable $Z \triangleq f_{\text{enc}}(X)$, and a neural network decoder $f_{\text{dec}}$ maps $Z$ back to a reconstruction $\hat{X} \triangleq f_{\text{dec}}(Z)$. A learned probability distribution $p_Z$, also known as the *entropy model*, is used to model the marginal distribution of $Z$. The learning objective is to minimize the rate-distortion (R-D) loss:

$$\min \mathbb{E}_{X \sim p_{\text{data}}} \left[ -\log_2 p_Z(Z) + \lambda \cdot d(X, \hat{X}) \right], \quad (1)$$

where $d$ is a distortion metric (*e.g.*, mean squared error), $\lambda$ is the Lagrange multiplier trading off between rate and distortion, and the minimization is over the network parameters of $f_{\text{enc}}$, $f_{\text{dec}}$, and $p_Z$. This framework has also been extended to *variable-rate compression* [12, 14], where the encoder, decoder, and entropy model are conditioned on $\lambda$. During variable-rate training, the model parameters are optimized for various $\lambda$ sampled from a distribution $p_\Lambda$:

$$\min \mathbb{E}_{X \sim p_{\text{data}}, \Lambda \sim p_\Lambda} \left[ -\log_2 p_Z(Z|\Lambda) + \Lambda \cdot d(X, \hat{X}) \right] \quad (2)$$

$$\text{where } Z \triangleq f_{\text{enc}}(X; \Lambda), \quad \hat{X} \triangleq f_{\text{dec}}(Z; \Lambda). \quad (3)$$

Existing research in LIC can be categorized into several groups. A major group focuses on designing expressive architectures for $f_{\text{enc}}$ and $f_{\text{dec}}$, such as convolutional and transformer-based ones [10, 11, 14, 20, 23, 27, 33, 37–39, 43, 52, 54, 59]. Another line of research lies in designing the entropy model $p_Z$, such as autoregressive [22, 40, 41, 44] and hierarchical models [3, 14, 15, 25]. Other research includes, *e.g.*, quantization methods [17, 18, 21, 55, 57] and variable-rate compression methods [6, 9, 12, 31, 48, 53]. To our knowledge, none of these existing methods are designed to work with continual learning as in our work.

The most related line of research to this paper is ***content-adaptive image compression***, where the goal is to adapt a compressor to new images or new target rates in a per-image fashion. Solutions to this problem include encoder-side optimization and decoder-side adaptation. Encoder-side optimization methods [7, 19, 56] directly optimize the R-D objective w.r.t. $Z$ during encoding. Decoder-side adaptation methods [42, 47, 50], on the other hand, include parameter-efficient neural network modules in the bitstream, which are executed on the decoder side to improve decoding. Among them, many methods require computationally expensive, iterative optimization during encoding.

The scope of this paper is distinct from content-adaptive image compression in several ways: (a) our goal is to incrementally train the compressor parameters in-place without introducing additional parameters, and (b) as opposed to per-image optimization during testing, our method employs a one-time training procedure and induces no additional computational cost at test time. Our research is complementary to content-adaptive image compression, and they can be combined to further improve the performance.

### 2.2. Knowledge replay in continual learning

Continual learning has been widely studied for image classification [13] and semantic segmentation [8], which aim to learn a sequence of new tasks incrementally without forgetting the previously learned knowledge. Among existing methods for continual learning, *replay-based* methods [34, 36, 45] have emerged as particularly effective,
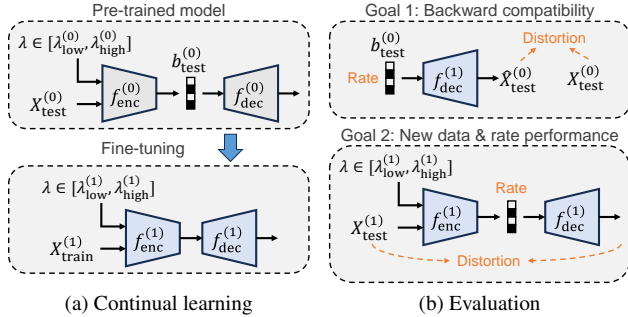
Figure 2. Problem definition. Fig. 2a shows the fine-tuning process of a pre-trained model, and Fig. 2b shows the two evaluation criteria: backward compatibility and new data/rate performance. Entropy models are kept frozen and omitted in the figures.

which applies an additional memory buffer to store exemplar data from learned tasks and then integrate with new task data to perform knowledge rehearsal during continual learning. Our proposed knowledge replay method adopts a similar principle. However, image compression contains unique challenges distinct from typical computer vision tasks, making existing continual learning strategies inapplicable to deploy in image compression. Thus, a tailored knowledge replay strategy is required for our problem scenario.

## 3. Problem Description

In this section, we formulate the problem (Sec. 3.1), discuss the backward compatibility issue (Sec. 3.2), and motivate the design of our method (Sec. 3.3).

### 3.1. Continual learning for compression

Assume that we have a pre-trained variable-rate model with supported $\lambda \in [\lambda_{\text{low}}^{(0)}, \lambda_{\text{high}}^{(0)}]$. We use $f_{\text{enc}}^{(0)}, f_{\text{dec}}^{(0)}, p_Z^{(0)}$ to denote the pre-trained model's encoder, decoder, and entropy model, respectively. Also, assume that we have used the pre-trained model to compress a test set of images $X_{\text{test}}^{(0)}$ and obtained the corresponding bitstreams $b_{\text{test}}^{(0)}$, as illustrated in Fig. 2a (top). This situation well-silumates a typical image storage application with a learned image compressor.

We now aim at fine-tuning the model with new data $X_{\text{train}}^{(1)}$ and a new rate range determined by $[\lambda_{\text{low}}^{(1)}, \lambda_{\text{high}}^{(1)}]$. Note that the new data and rate range may or may not be the same as the old ones. Similarly, let $f_{\text{enc}}^{(1)}, f_{\text{dec}}^{(1)}, p_Z^{(1)}$ denote the new model components, as illustrated in Fig. 2a (bottom). We expect the new model to achieve two goals:

1. **Backward compatibility:** The new model should be capable of decoding the bitstreams produced by the old model (Fig. 2b, top).
2. **New-data (or new-rate) performance**: The new model should perform better than the old one on new data and new rates (Fig. 2b, bottom).
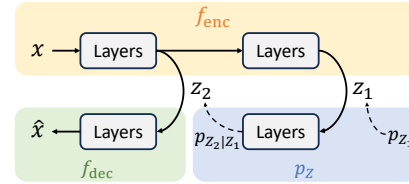


Figure 3. The Hyperprior model architecture [3], where the layers are categorized into three groups: $f_{\text{enc}}$, $f_{\text{dec}}$, and $p_Z$.

| | **Number of parameters** | | | |
|---|---|---|---|---|
| | Total | $f_{\text{enc}}$ | $f_{\text{dec}}$ | $p_Z$ |
| M-S Hyp. [41] | 17.6M | 5.9M | 2.5M | 8.2M (47%) |
| GMA [11] | 26.6M | 5.8M | 11.7M | 9.0M (34%) |
| ELIC [23] | 33.8M | 9.7M | 7.3M | 16.7M (49%) |
| STF [59] | 99.9M | 11.3M | 7.1M | 81.4M (81%) |
| TCM [33] | 45.2M | 3.3M | 6.8M | 35.2M (78%) |
| MLIC++ [26] | 116.7M | 6.3M | 12.0M | 98.4M (84%) |
| Our model | 35.5M | 18.7M | 11.9M | **4.9M (14%)** |

Table 1. Many existing methods employ a parameter-expensive entropy model. We propose an architecture with a lightweight entropy model, which makes continual learning more effective (since more parameters are learnable in the fine-tuning phase).

### 3.2. Backward compatibility of entropy decoding

As mentioned in Fig. 1, fine-tuning the model end-to-end breaks the backward compatibility of the decoder. The primary reason lies in entropy coding: range-based entropy coding algorithms [16, 46], which are commonly used in modern neural compressors, are known to be sensitive to the probability distribution of the encoded symbols. As the new entropy model $p_Z^{(1)}$ is different from the old one $p_Z^{(0)}$, the new model is not able to perform entropy coding correctly using the old bitstreams $b_{\text{test}}^{(0)}$, and thus the correct (quantized) latent variables $Z_{\text{test}}^{(0)}$ cannot be obtained. In fact, recent research [2, 29, 49] has shown that even a small change (*e.g.*, a floating point round-off error) in the entropy model can lead to failure in entropy decoding.

### 3.3. Freezing the entropy model during fine-tuning

To avoid the aforementioned issue, the entropy model $p_Z$ needs to be kept unchanged throughout fine-tuning. Then, the latent variables $Z_{\text{test}}^{(0)}$ are guaranteed to be recovered losslessly from the old bitstreams, and the problem reduces to continually learning the decoder $f_{\text{dec}}$ without forgetting the old knowledge (*i.e.*, decoding $Z_{\text{test}}^{(0)}$). With our proposed training strategy (Sec. 4.1), we show that this can be achieved for many existing neural compressors.

We also notice that the entropy model $p_Z$ is often the largest component in many existing compressors, most of which are based on the Hyperprior model [3, 41] (Fig. 3). We found that their entropy model takes up a significant portion (*e.g.*, 84% for MLIC++[26]) of the model parameters, as shown in Tab. 1. Consequently, a large proportion of model parameters need to be fixed during continual learn-

ing, which may impair the new-data (or new-rate) performance. Motivated by this, we design a model architecture (Sec. 4.2) that employs a lightweight entropy model (Tab. 1, last row), which makes continual learning more effective.

# 4. Method

As just mentioned, our methods include two independent components: the knowledge replay-based training strategy (Sec. 4.1) and a neural network architecture specifically designed for continual learning (Sec. 4.2). We now present them sequentially in detail.

## 4.1. Continual learning with knowledge replay

Following the notation introduced in Sec. 3.1, we denote a pre-trained model as $f_{\text{enc}}^{(0)}, f_{\text{dec}}^{(0)}, p_Z$, which are trained on data $X_{\text{train}}^{(0)}$ with $\lambda \in [\lambda_{\text{low}}^{(0)}, \lambda_{\text{high}}^{(0)}]$. We drop the superscript for $p_Z$ since it is kept frozen throughout fine-tuning. Inspired by the idea of knowledge rehearsal with exemplars in class-incremental learning methods [45], we use the old training data $X_{\text{train}}^{(0)}$ and the old encoder $f_{\text{enc}}^{(0)}$ to perform "*knowledge replay*" in the fine-tuning process. To this end, we store $f_{\text{enc}}^{(0)}$ along with $X_{\text{train}}^{(0)}$ within a dedicated memory buffer before the fine-tuning stage. Note that we do not pose restrictive assumptions on training resources and allow access to the entire training set $X_{\text{train}}^{(0)}$ during fine-tuning.

In the fine-tuning process, our knowledge replay-based training objective contains two terms. The first term, $\ell_{\text{new}}$, is the standard R-D loss for the new training data $X_{\text{train}}^{(1)}$ with the new $\lambda$ value range $[\lambda_{\text{low}}^{(1)}, \lambda_{\text{high}}^{(1)}]$:

$$\ell_{\text{new}} \triangleq \mathbb{E}\left[ R^{(1)} + \Lambda^{(1)} \cdot D^{(1)} \right], \quad \text{where} \quad (4)$$

$$R^{(1)} \triangleq -\log_2 p_Z(Z^{(1)}|\Lambda^{(1)}) \quad (5)$$

$$D^{(1)} \triangleq d(X_{\text{train}}^{(1)}, \hat{X}_{\text{train}}^{(1)}). \quad (6)$$

In (4), the expectation is w.r.t. $X_{\text{train}}^{(1)}$ and $\Lambda^{(1)}$, where $\Lambda^{(1)}$ is a random variable with the support being $[\lambda_{\text{low}}^{(1)}, \lambda_{\text{high}}^{(1)}]$. Its probability density, $p_\Lambda^{(1)}$, controls how $\lambda$ is sampled during training. Intuitively, minimizing $\ell_{\text{new}}$ adapts the model parameters to new data and new rates, but it is not sufficient to maintain backward compatibility.

The other term in our loss function encourages backward compatibility of the model parameters through knowledge replay of the old data $X_{\text{train}}^{(0)}$ and the old encoder $f_{\text{enc}}^{(0)}$. Specifically, we use $f_{\text{enc}}^{(0)}$ to encode $X_{\text{train}}^{(0)}$, which gives the corresponding (quantized) latent variables $Z_{\text{train}}^{(0)}$. Then, the current decoder $f_{\text{dec}}^{(1)}$ decodes $Z_{\text{train}}^{(0)}$, and the reconstruction $\hat{X}_{\text{train}}^{(0)}$ is compared with $X_{\text{train}}^{(0)}$ to compute the knowledge re-
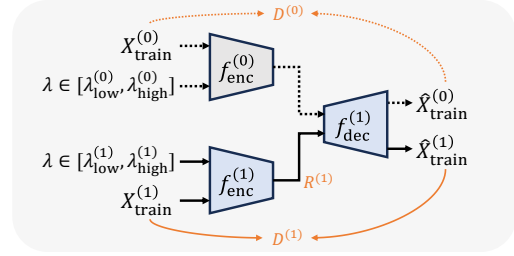


Figure 4. Our knowledge replay-based training strategy contains two components: a distortion loss that encourages backward compatibility (top, dashed lines), and the standard R-D loss for new data and new rates (bottom, solid lines). The entropy model $p_Z$ is kept frozen and is omitted in the figure.

play loss function, $\ell_{\text{KR}}$, defined as:

$$\ell_{\text{KR}} \triangleq \mathbb{E}\left[ \Lambda^{(0)} \cdot D^{(0)} \right], \quad \text{where} \quad (7)$$

$$D^{(0)} \triangleq d\left( X_{\text{train}}^{(0)}, f_{\text{dec}}^{(1)}\left( f_{\text{enc}}^{(0)}(X_{\text{train}}^{(0)}) \right) \right). \quad (8)$$

The expectation is w.r.t. $X_{\text{train}}^{(0)}$ and $\Lambda^{(0)}$, where $\Lambda^{(0)} \sim p_\Lambda^{(0)}$ controls how $\lambda$ is sampled during knowledge replay. In our experiments, we choose $p_\Lambda^{(0)}$ to be the same as the one used in pre-training. Note that there is no rate term in $\ell_{\text{KR}}$ since the replayed encoder $f_{\text{enc}}^{(0)}$ is fixed, and thus the rate term is constant w.r.t. the model parameters being trained. By replaying the old data and the old encoder, the decoder network effectively retains backward compatibility, as shown in our experiments (Sec. 5.3).

Fig. 4 illustrates our knowledge replay strategy for one training iteration. Our training objective is to minimize a weighted summation of the two terms, $\ell_{\text{new}}$ and $\ell_{\text{KR}}$, with a scalar hyperparameter $\alpha \in [0, 1]$:

$$\min (1 - \alpha) \cdot \ell_{\text{new}} + \alpha \cdot \ell_{\text{KR}}. \quad (9)$$

The hyperparameter $\alpha$ controls the weight of replayed data in the training objective and can be tuned to achieve a desired trade-off between backward compatibility and new data/rate performance (shown in Sec. 5.4). Our knowledge replay strategy is general and can be applied to various model architectures, enabling backward-compatible continual learning of those models (shown in Sec. 5.3).

## 4.2. Model architecture

Since freezing the entropy model is necessary for backward compatibility, we propose a model architecture that employs a lightweight entropy model by design so that most parameters in the model can stay learnable. An overview of the architecture is shown in Fig. 5. Our model is inspired by the hierarchical residual coding architecture [14, 18, 58], but we decouple the entropy model ($p_Z$) and the decoder

($f_{dec}$) into two separate branches, resulting in a reduced entropy model size. We now describe each component of the model in detail.

**Encoding:** The encoding process involves a bottom-up pass through the encoder $f_{enc}$ and a top-down pass through the entropy model $p_Z$. Given an input image $x$, encoding begins with $f_{enc}$, a neural network consisting of a sequence of downsampling and residual layers that produces a hierarchy of image features (denoted as $h_i$ in the figure). Specifically, for an input image with $256 \times 256$ pixels, $f_{enc}$ produces $N = 4$ features with spatial dimentions $32 \times 32$, $16 \times 16$, $8 \times 8$, and $4 \times 4$. All layers are convolutional, so the spatial dimensions scale accordingly for images of different sizes. Then, the entropy model $p_Z$ starts with a constant $e_0$ and iteratively updates it using the features $h_i$ from $f_{enc}$. In each stage, $e_{i-1}$ is upsampled to the same spatial dimension as $h_i$ and concatenated with $h_i$. The concatenated features are then fed into a sequence of layers to produce $z_i$, the latent variable (which will be entropy coded) for the $i$-th stage. $z_i$ is then aggregated with the upsampled $e_{i-1}$ through a linear layer and addition operation, and the result is denoted as $e_i$ and passed to the next stage. Note that in each stage, the entropy model also estimates the probability distribution of $z_i$ given $z_{<i}$ (with $z_{<i} \triangleq \{z_1, ..., z_{i-1}\}$), which is used for entropy coding.

**Probabilistic model, quantization, and entropy coding** are performed in the same way as for the discretized Gaussian model in Hyperprior-based methods [3, 41]. As a brief recap, the entropy model predicts a mean $\mu_i$ and a scale $\sigma_i$ for each latent variable $z_i$, and the probability model for $z_i$ is a discretized Gaussian distribution:

$$p_i(z_i) \triangleq p_{Z_i|Z_{<i}}(z_i \mid z_{<i}) \tag{10}$$

$$= \int_{z_i-0.5}^{z_i+0.5} \mathcal{N}(t; \mu_i, \sigma_i^2) \, dt, \tag{11}$$

where the dependence on $z_{<i}$ is through $\mu_i$ and $\sigma_i$. During testing, the residual between $z_i$ and $\mu_i$ is quantized to the nearest integer, and during training, quantization is simulated by additive uniform noise. Each stage $i$ produces a separate bitstream corresponding to $z_i$ (using the asymmetric numeral systems [16] algorithm), and all stages are executed sequentially for $i = 1, ..., N$.

**Decoding:** Given encoded bitstreams, the decoding process mirrors the encoding process in reverse. Firstly, the entropy model is executed top-down to iteratively predict $p_{Z_i|Z_{<i}}$, based on which the bitstreams are entropy-decoded to obtain $z_i$. Then, the decoder $f_{dec}$ is executed top-down to reconstruct the image. Starting with a constant $r_0$, the decoder iteratively updates $r_i$ using $z_i$ and $e_i$ in each stage, as shown in the right pane of Fig. 5. In each stage and the final layer, residual layers and upsampling layers are applied to restore the image to its original resolution.
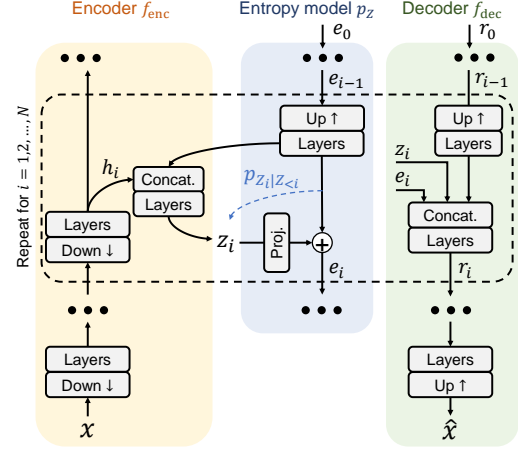


Figure 5. Our model adopts the hierarchical residual coding architecture [14, 18, 58] but decouples the entropy model ($p_Z$) and the decoder ($f_{dec}$) into two branches. *Up* ↑ denotes upsampling, *Down* ↓ denotes down-sampling, and *Proj.* denotes a linear projection layer. Detailed layer configurations are in Appendix, Sec. 7
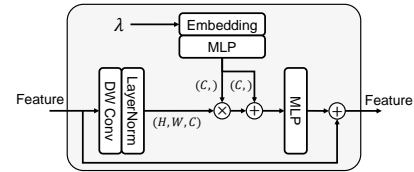


Figure 6. Each layer in our model is a ConvNeXt block [35] conditioned on $\lambda$ through an affine transformation. In the figure, $(H, W, C)$ denotes height, width, and channel dimensions.

**Rate-conditional network layers:** Fig. 6 shows the details of each layer (*i.e.*, residual block) in our model. Our model employs the ConvNeXt module [35] as the basic building blocks. To achieve variable-rate compression, we adopt the conditional convolution technique [12, 14], which applies an adaptive affine transformation to the convolutional layer output (after layer normalization) to control the rate based on the input $\lambda$.

**Variable-rate training:** The (pre-)training objective is to minimize the standard R-D loss for variable rate compression (Eq. (2)), except that the rate term consists of the sum of the rates for all latent variables:

$$\min \mathbb{E}_{X,\Lambda} \left[ \sum_{i=1}^{N} -\log_2 p_i(Z_i|\Lambda) + \Lambda \cdot d(X, \hat{X}) \right], \tag{12}$$

where $X$ follows the training data distribution, and $\Lambda$ follows $p_\Lambda$, a continuous probability distribution that controls the sampling strategy of $\lambda$ during training. After the pre-training phase, we freeze $p_Z$ and apply our knowledge replay training strategy presented in Sec. 4.1 to achieve backward-compatible fine-tuning.

# 5. Experiments

Our experiments compare various fine-tuning strategies as well as various model architectures for continual learning of image compression. We begin by describing the setup (Sec. 5.1) and baseline methods (Sec. 5.2). Then, Sec. 5.3 presents the main experimental results for our proposed methods. Finally, we provide additional experiments to analyze the effectiveness of knowledge replay (Sec. 5.4) and our model architecture (Sec. 5.5).

## 5.1. Experiment setup

We consider two continual learning scenarios in image compression: data-incremental learning and rate-incremental learning. In the former, pre-trained models are fine-tuned on a new dataset, and in the latter, models are fine-tuned with a larger rate range (either going higher or lower). Detailed configurations are shown in Tab. 2, and the datasets and metrics used are listed below.

**Datasets:** we use the COCO [32] dataset *train2017* split for pre-training all models. The dataset contains 118,287 images with around $640 \times 420$ pixels. We randomly crop the images to $256 \times 256$ patches during training. For data-incremental learning, we adopt the CelebA-HQ dataset [30] at $256 \times 256$ pixels, a commonly-used human face image dataset for generative image modeling [24]. The dataset consists of 30,000 images, where 24,000 are for training, 3,000 for validation, and the remaining 3,000 for testing.

**Metrics:** We use bits per pixel (bpp), peak signal-to-noise ratio (PSNR, computed for the RGB space), and BD-Rate [5] to measure compression performance, all of which are standard metrics for image compression. As described in Sec. 3, we evaluate each method for two objectives:

- **Backward compatibility:** We use the fine-tuned model to decode $b_{\text{test}}^{(0)}$, the bitstreams encoded by the pre-trained model, to obtain reconstructions $\hat{X}_{\text{test}}^{(0)}$. The bpp is computed for $b_{\text{test}}^{(0)}$ (which is a constant independent of fine-tuning strategies), and the PSNR is computed between the reconstructions $\hat{X}_{\text{test}}^{(0)}$ and the original images $X_{\text{test}}^{(0)}$.
- **New data & rate performance:** We compress the new data $X_{\text{test}}^{(1)}$ with the new rates, determined by the $\lambda$ value range $[\lambda_{\text{low}}^{(1)}, \lambda_{\text{high}}^{(1)}]$, to obtain bpp and PSNR metrics.

We report all results in terms of BD-Rate in the main paper (due to space constraints), and we provide the corresponding PSNR-bpp curves in the Appendix.

## 5.2. Methods in comparison

In addition to our model, we choose *Mean & Scale Hyperprior* [41] (MSH) and *Gaussian Mixture & Attention* [11] (GMA) as two base models for our experiments. These two are commonly used and representative models for learned image compression, and we believe the experimental conclusions based on them can be generalized to other existing

|  | $[\lambda_{\text{low}}^{(0)}, \lambda_{\text{high}}^{(0)}]$ | $X_{\text{train}}^{(0)}$ | $X_{\text{test}}^{(0)}$ |
|---|---|---|---|
| Pre-training | [32, 1024] | COCO | Kodak |

|  | $[\lambda_{\text{low}}^{(1)}, \lambda_{\text{high}}^{(1)}]$ | $X_{\text{train}}^{(1)}$ | $X_{\text{test}}^{(1)}$ |
|---|---|---|---|
| Data-incremental | [32, 1024] | CelebA | CelebA |
| Rate-incremental (low → high) | [32, 4096] | COCO | Kodak |
| Rate-incremental (high → low) | [4, 1024] | COCO | Kodak |

Table 2. Experiment configurations. We start with a pre-trained model (*pre-training*) and fine-tune it either with new data (*data-incremental*) or new rates (*rate-incremental*).

|  | BD-Rate (%) w.r.t. VTM 22.0 ↓ | | |
|---|---|---|---|
|  | Old bitstreams (Kodak) | New data (CelebA) | Avg. |
| MSH-VR, pre-trained | 26.7 | 17.1 | 21.9 |
| MSH-VR w/ FT Enc. | 26.7 | 14.8 | 20.8 |
| MSH-VR w/ FT Enc. & Dec. | 419.1 | **8.95** | 214.0 |
| MSH-VR w/ KR (ours) | **20.1** | 9.57 | **14.8** |
| GMA-VR, pre-trained | 4.43 | -9.23 | -2.40 |
| GMA-VR w/ FT Enc. | 4.43 | -10.5 | -3.04 |
| GMA-VR w/ FT Enc. & Dec. | 354.6 | **-17.3** | 168.7 |
| GMA-VR w/ KR (ours) | **4.33** | -14.4 | **-5.04** |
| Our model, pre-trained | 1.87 | -13.2 | -5.67 |
| Our model w/ FT Enc. | 1.87 | -14.6 | -6.37 |
| Our model w/ FT Enc. & Dec. | 262.9 | **-19.0** | 122.0 |
| Our model w/ KR | **0.87** | -16.6 | **-7.87** |

Table 3. Data-incremental learning (COCO → CelebA) results. PSNR-bpp curves are provided in Appendix, Fig. 8.

models. Since variable-rate compression is required to perform rate-incremental learning, we construct a variable-rate version for each of them and train it in the same way as for our model. The resulting models are referred to as *MSH-VR* and *GMA-VR*, respectively. Sec. 8 in the Appendix provides pre-training and fine-tuning hyperparameters, and Sec. 9 in the Appendix provides details on the variable-rate baselines compared to their fixed-rate models.

We apply the following fine-tuning strategies for each model and compare their performance:

- **Pre-traiend model**: Using the pre-trained model without fine-tuning with new data or rates is the simplest baseline.
- **Fine-tuning the encoder only (FT Enc.)**: We fine-tune the encoder with new data while keeping other parameters frozen. Since the entropy model and the decoder are never changed, backward compatibility is guaranteed.
- **Fine-tuning both the encoder and decoder (FT Enc. & Dec.)**: We fine-tune all model parameters except for the entropy model parameters. Since the decoder changes, backward compatibility may be lost.
- **Our approach: knowledge replay (KR).** We fine-tune the model's encoder and decoder with the proposed knowledge replay (KR) strategy applied.

## 5.3. Experimental results

**Data-incremental learning**. Tab. 3 show the results (pre-trained on COCO, fine-tuned on CelebA). We begin by

| | Kodak BD-Rate (%) w.r.t. VTM 22.0 ↓ | | |
|---|---|---|---|
| | Old bitstreams: bpp ≈ (0.1,0.9) | New rate: bpp ≈ (0.1,1.6) | Avg. |
| MSH-VR, pre-trained | 26.7 | - | - |
| MSH-VR w/ FT Enc. & Dec. | 282.2 | 23.4 | 152.80 |
| MSH-VR w/ KR (ours) | **19.6** | **17.3** | **18.45** |
| GMA-VR, pre-trained | 4.43 | - | - |
| GMA-VR w/ FT Enc. & Dec. | 64.1 | 4.56 | 34.33 |
| GMA-VR w/ KR (ours) | **3.39** | **2.34** | **2.87** |
| Our model, pre-trained | 1.87 | - | - |
| Our model w/ FT Enc. & Dec. | 17.7 | 1.45 | 9.58 |
| Our model w/ KR | **0.96** | **0.70** | **0.83** |

Table 4. Rate-incremental learning (low → high) results. PSNR-bpp curves are provided in Appendix, Fig. 9.

| | Kodak BD-Rate (%) w.r.t. VTM 22.0 ↓ | | |
|---|---|---|---|
| | Old bitstreams: bpp ≈ (0.1,0.9) | New rate: bpp ≈ (0.03,0.9) | Avg. |
| MSH-VR, pre-trained | 26.7 | - | - |
| MSH-VR w/ FT Enc. & Dec. | 35.1 | 37.6 | 36.4 |
| MSH-VR w/ KR (ours) | **18.9** | **29.9** | **24.4** |
| GMA-VR, pre-trained | 4.43 | - | - |
| GMA-VR w/ FT Enc. & Dec. | 10.0 | 9.11 | 9.56 |
| GMA-VR w/ KR (ours) | **1.75** | **6.92** | **4.34** |
| Our model, pre-trained | 1.87 | - | - |
| Our model w/ FT Enc. & Dec. | 14.33 | 5.18 | 9.76 |
| Our model w/ KR | **0.86** | **4.26** | **2.56** |

Table 5. Rate-incremental learning (high → low) results. PSNR-bpp curves are provided in Appendix, Fig. 10.

| | Fine-tuning | | BD-Rate w.r.t. VTM 22.0 | | |
|---|---|---|---|---|---|
| Config. | Data | KR loss | Old bits. | New data | Avg. |
| 0 | - | - | 1.87 | -13.2 | -5.67 |
| 1 | CA | | 262.9 | **-19.0** | 121.9 |
| 2 | CA + COCO | | 23 | -16.9 | 3.05 |
| 3 | CA | ✓ | 3.67 | -16.3 | -6.32 |
| 4 (ours) | CA + COCO | ✓ | **0.87** | -16.6 | **-7.87** |

Table 6. Ablative analysis of our knowledge replay-based training strategy for data-incremental learning (COCO → CelebA). For the *"data"* column, *CA* denotes the CelebA dataset.

comparing the fine-tuning strategy for the *MSH-VR* model. Firstly, fine-tuning the encoder (*FT Enc.*) does not provide a significant improvement for new data BD-Rate (17.1% → 14.8%), which is expected since fine-tuning the encoder only reduces the amortization gap [56] and does not improve the model's capacity. When fine-tuning both the encoder and decoder (*FT Enc. & Dec.*), the new data BD-Rate is improved by a much larger margin (17.1% → 8.95%), indicating the importance of updating the decoder. However, this came at the cost of backward compatibility, as shown by the significant increase in BD-Rate for old bitstreams (26.7% → 419.1%). This indicates that, while the decoder fits the new data well, it becomes incompatible with the old bitstreams. With our knowledge replay strategy (*MSH-VR w/ KR*), we are able to achieve a competitive new data performance (9.57% BD-Rate) without sacrificing backward compatibility. Notably, fine-tuning with our strategy also improves the performance on old bitstreams, which is not the case for the other strategies. On average, the knowledge replay strategy clearly outperforms the other ones. These observations are consistent with the results for *GMA-VR* and our model, demonstrating the effectiveness of knowledge replay in data-incremental learning. Also, when comparing different model architectures, our model achieves the best performance overall in terms of all metrics.

**Rate-incremental learning**. Tab. 4 presents the re-

sults for rate-incremental learning (from low rates to higher rates). In rate-incremental experiments, we omit the *FT Enc.* baseline, since fine-tuning the encoder alone cannot effectively extend the rate range of any considered models (see Appendix, Sec. 10.2 for details). Starting with the *MSH-VR* model, we observe that fine-tuning both the encoder and decoder (*FT Enc. & Dec.*) is able to extend the operational rate range of the model with a similar BD-Rate w.r.t. VTM for the new rates. However, the performance on old bitstreams is significantly degraded, similar to the observation in data-incremental learning experiments. By applying our knowledge replay strategy (*MSH-VR w/ KR*), in contrast, the model is able to achieve a competitive BD-Rate for the new rates (17.3%) while maintaining backward compatibility. Again, the results for *GMA-VR* and our model show a similar pattern. Overall, our model with KR outperforms the baselines, validating its effectiveness in rate-incremental learning. For rate-incremental learning from high rates to low rates (Tab. 5), the above observations stay the same.

### 5.4. Experimental analysis: knowledge replay

The effectiveness of the proposed knowledge replay strategy has already been verified in the previous experiments. We now provide additional experiments to answer the following questions that aim to analyze the individual components of our knowledge replay strategy.

**What contributes to the backward compatibility?** In data-incremental learning, there are two components in our knowledge replay strategy that may help backward compatibility: the replayed training data and the knowledge replay loss. To analyze the contribution of each component, we freeze the entropy model parameters of our model and fine-tune it with the two components separately. To analyze the contribution of each component, we start from "*Our model w/ FT Enc. & Dec.*" and apply the two components one by one. Tab. 6 shows the results. *Config. 0* is the pre-trained model, and *Config. 1* is the "*FT Enc. & Dec.*" baseline that does not retain backward compatibility. With replayed data (*Config. 2*), backward compatibility is largely improved (262.9% → 23% BD-Rate) but is still worse than the pre-trained model. When the knowledge replay loss is applied (*Config. 3* and *Config. 4*), the performance on old bitstreams becomes comparable to the pre-trained model.

| $\alpha$ | 0.0 | 0.25 | 0.5 | 0.75 | 1.0 |
|---|---|---|---|---|---|
| BD-rate: old bitstreams | 262.9 | 1.40 | 0.87 | 0.72 | **0.63** |
| BD-rate: new data | **-19.0** | -16.8 | -16.6 | -16.3 | -14.4 |
| Avg. BD-rate | 122.0 | -7.7 | **-7.9** | -7.8 | -6.9 |

Table 7. Data-incremental learning (COCO → CelebA) results for our model with varying $\alpha$. The scalar $\alpha \in [0, 1]$ controls the ratio of replayed data in fine-tuning, where $\alpha = 0.0$ means no replay, and $\alpha = 1.0$ means no new data. BD-rate is w.r.t. VTM 22.0.

| | | Latency (in seconds) | | | | | |
|---|---|---|---|---|---|---|---|
| | Params. | Entropy Coding | | Network (CPU) | | Network (GPU) | |
| | | Enc. | Dec. | Enc. | Dec. | Enc. | Dec. |
| MSH + VR | 19.2M | 0.026 | 0.068 | 0.318 | 0.325 | 0.004 | 0.010 |
| GMA + VR | 33.4M | 3.072 | 6.302 | 0.941 | 1.221 | 0.006 | 0.022 |
| Ours | 35.5M | 0.046 | 0.052 | 0.474 | 0.393 | 0.026 | 0.011 |

*Hardware: Intel 10700K CPU (using four threads) and Nvidia Quadro 6000 GPU.
*Latency is the average time to encode/decode a Kodak image (768×512 pixels), averaged over all 24 images. Time includes entropy coding.

Table 8. Computational complexity of the model architectures used in our experiments.

| | BD-Rate (%) w.r.t. VTM 22.0 ↓ | | |
|---|---|---|---|
| | Old bitstreams (Kodak) | New data (CelebA) | Avg. |
| Sequential $p_Z$ and $f_{\text{dec}}$ | 4.42 | -12.8 | -4.19 |
| Sequential $p_Z$ and $f_{\text{dec}}$ w/ KR | 3.18 | -15.8 | -6.31 |
| Parallel $p_Z$ and $f_{\text{dec}}$ | 1.87 | -13.2 | -5.67 |
| Parallel $p_Z$ and $f_{\text{dec}}$ w/ KR | **0.87** | **-16.6** | **-7.87** |

Table 9. Comparing architecture variants of our model in terms of data-incremental learning (COCO → CelebA).

to the increase in computational complexity.

**Impact of decoupling the entropy model and the decoder.** Recall that in order to reduce the parameters of $p_Z$, our architecture decouples $p_Z$ and $f_{\text{dec}}$ into two parallel branches. To analyze the impact of doing so, we construct a variant of our model where $p_Z$ and $f_{\text{dec}}$ are executed in sequential (like in Hyperprior-based models; see Fig. 3). For a fair comparison, the sequential version has the same number of latent feature channels and the same number of parameters as the parallel model, and training is performed in the same way as for all previous experiments. We show the data-incremental learning results in Tab. 9. We observe that our model (*Parallel $p_Z$ and $f_{dec}$ w/ KR*) achieves a better performance than the sequential version on both old bitstreams and new data, which verifies our design.

## 5.6. Discussion

Our experiments show that neural image compressors can adapt to new data and rates in a backward-compatible manner by using the proposed training strategy. In addition to continual learning applications, this observation offers insights into related research such as the standardization of learned image compression. Despite recent attempts toward this goal [1, 39], it remains an open question about which part of a selected neural compressor needs to be standardized. Our findings suggest a possible direction: only the entropy model (instead of the entire model architecture and parameters) needs to be standardized, and other components (*e.g.*, the decoder network) could be fine-tuned over time with backward-compatible training strategies.

## 6. Conclusion

This paper presents two approaches: a knowledge replay-based training strategy and a neural network architecture, for continual learning of image compression. Our knowledge replay strategy enables existing compressors to adapt to new data and target rates while ensuring that previously compressed bitstreams remain decodable. Through extensive experiments, we conclusively answer the question raised at the beginning of the paper: neural network-based image compressors can be learned continually in a backward-compatible manner, achieving improved performance on new data, new rates, and old bitstreams.

We conclude that both the replayed data and the loss function contribute to backward compatibility, and the knowledge replay loss is more important among the two.

**What is the impact of $\alpha$ in the knowledge replay loss function (Eq. (9))?** We train our model for data-incremental learning with varying $\alpha$, the scalar that controls the ratio of replayed data in each training iteration. Results are shown in Tab. 7. Firstly, it is clear that as $\alpha$ grows from 0 to 1, the performance on old bitstreams monotonically improves (*i.e.*, the BD-Rate decreases). Recall that when $\alpha = 0$, no replay is performed, and the model is trained with only the new data; when $\alpha = 1$, the model is trained with only the replayed data. The results are thus consistent with the intuition that more replay leads to better backward compatibility. When it comes to the new data performance, the trend is reversed: as $\alpha$ grows, the performance on new data monotonically degrades, which again matches the intuition. On average, the performance is comparable for $\alpha \in [0.25, 0.75]$. We conclude that our approach is insensitive to the choice of $\alpha$ and can achieve a good trade-off between backward compatibility and new data performance. Our experiments choose $\alpha = 0.5$ by default, but in practice, the choice of $\alpha$ can be treated as a hyperparameter and determined by the application requirements.

## 5.5. Experimental analysis: model architecture

**Computational complexity.** Tab. 8 shows the computational complexity of our model and the two baselines. The metrics include the number of parameters, entropy coding latency, and neural network forward pass latency. Except for a few exceptions (the parameter count of *MSH-VR*, the entropy coding latency of *GMA-VR*, and the GPU encoding latency of our model), all methods are mostly comparable in terms of computational complexity. Thus, we conclude that the performance improvement of our model is not due

**Limitations and future work.** Our work serves as a preliminary study on continual learning for image compression. Despite its effectiveness, our knowledge replay strategy assumes unconstrained training resources, which may not be true in practice. Also, we focus on the decoder's backward compatibility, while the same problem can be studied for the encoder. For future work, a possible direction is to extend our two-step method (pre-training and fine-tuning) to multi-step continual learning scenarios.

# References

[1] João Ascenso, Elena Alshina, and Touradj Ebrahimi. The jpeg ai standard: Providing efficient human and machine visual data consumption. *IEEE MultiMedia*, 30(1):100–111, 2023. 8

[2] Johannes Ballé, Nick Johnston, and David Minnen. Integer networks for data compression with latent-variable models. *International Conference on Learning Representations*, 2018. 3

[3] J. Ballé, D. Minnen, S. Singh, S. Hwang, and N. Johnston. Variational image compression with a scale hyperprior. *International Conference on Learning Representations*, 2018. 2, 3, 5

[4] J. Ballé, P. A. Chou, D. Minnen, S. Singh, N. Johnston, E. Agustsson, S. Hwang, and Toderici G. Nonlinear transform coding. *IEEE Journal of Selected Topics in Signal Processing*, 15(2):339–353, 2021. 2

[5] Gisle Bjontegaard. Calculation of average psnr differences between rd-curves. *Video Coding Experts Group - M33*, 2001. 6

[6] Shilv Cai, Zhijun Zhang, Liqun Chen, Luxin Yan, Sheng Zhong, and Xu Zou. High-fidelity variable-rate image compression via invertible activation transformation. *Proceedings of the ACM International Conference on Multimedia*, pages 2021–2031, 2022. 2

[7] Joaquim Campos, Simon Meierhans, Abdelaziz Djelouah, and Christopher Schroers. Content adaptive optimization for neural image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019. 2

[8] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Bulo, Elisa Ricci, and Barbara Caputo. Modeling the background for incremental learning in semantic segmentation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9233–9242, 2020. 2

[9] Tong Chen and Zhan Ma. Variable bitrate image compression with quality scaling factors. *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 2163–2167, 2020. 2

[10] T. Chen, H. Liu, Z. Ma, Q. Shen, X. Cao, and Y. Wang. End-to-end learnt image compression via non-local attention optimization and improved context modeling. *IEEE Transactions on Image Processing*, 30:3179–3191, 2021. 2

[11] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7936–7945, 2020. 2, 3, 6, 4

[12] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Variable rate deep image compression with a conditional autoencoder. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3146–3154, 2019. 2, 5

[13] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3366–3385, 2022. 1, 2

[14] Zhihao Duan, Ming Lu, Jack Ma, Yuning Huang, Zhan Ma, and Fengqing Zhu. Qarv: Quantization-aware resnet vae for lossy image compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–15, 2023. 2, 4, 5

[15] Zhihao Duan, Ming Lu, Zhan Ma, and Fengqing Zhu. Lossy image compression with quantized hierarchical vaes. *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 198–207, 2023. 2

[16] Jarek Duda, Khalid Tahboub, Neeraj J. Gadgil, and Edward J. Delp. The use of asymmetric numeral systems as an accurate replacement for huffman coding. *Picture Coding Symposium*, pages 65–69, 2015. 3, 5

[17] Alaaeldin El-Nouby, Matthew J. Muckley, Karen Ullrich, Ivan Laptev, Jakob Verbeek, and Herve Jegou. Image compression with product quantized masked image modeling. *Transactions on Machine Learning Research*, 2023. 2

[18] Runsen Feng, Zongyu Guo, Weiping Li, and Zhibo Chen. Nvtc: Nonlinear vector transform coding. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6101–6110, 2023. 2, 4, 5

[19] Chenjian Gao, Tongda Xu, Dailan He, Yan Wang, and Hongwei Qin. Flexible neural image compression via code editing. *Advances in Neural Information Processing Systems*, 35:12184–12196, 2022. 2

[20] G. Gao, P. You, R. Pan, S. Han, Y. Zhang, Y. Dai, and H. Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021. 2

[21] Zongyu Guo, Zhizheng Zhang, Runsen Feng, and Zhibo Chen. Soft then hard: Rethinking the quantization in neural image compression. *Proceedings of the International Conference on Machine Learning*, 139:3920–3929, 2021. 2

[22] Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. Checkerboard context model for efficient learned image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14766–14775, 2021. 2

[23] Dailan He, Ziming Yang, Weikun Peng, Rui Ma, Hongwei Qin, and Yan Wang. Elic: Efficient learned image compression with unevenly grouped space-channel contextual adaptive coding. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5708–5717, 2022. 2, 3

[24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 6

[25] Yueyu Hu, Wenhan Yang, Zhan Ma, and Jiaying Liu. Learning end-to-end lossy image compression: A benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4194–4211, 2022. 2

[26] Wei Jiang and Ronggang Wang. MLIC$^{++}$: Linear complexity multi-reference entropy modeling for learned image compression. *ICML Workshop Neural Compression: From Information Theory to Applications*, 2023. 3

[27] Wei Jiang, Jiayu Yang, Yongqi Zhai, Peirong Ning, Feng Gao, and Ronggang Wang. Mlic: Multi-reference entropy model for learned image compression. *Proceedings of the ACM International Conference on Multimedia*, pages 7618–7627, 2023. 2

[28] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. 1

[29] Esin Koyuncu, Timofey Solovyev, Elena Alshina, and André Kaup. Device interoperability for learned image compression with weights and activations quantization. *Picture Coding Symposium*, pages 151–155, 2022. 3

[30] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5548–5557, 2020. 6, 1

[31] Jooyoung Lee, Seyoon Jeong, and Munchurl Kim. Selective compression learning of latent representations for variable-rate image compression. *Advances in Neural Information Processing Systems*, 35:13146–13157, 2022. 2

[32] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. *Proceedings of the European Conference on Computer Vision*, pages 740–755, 2014. 6, 1

[33] Jinming Liu, Heming Sun, and Jiro Katto. Learned image compression with mixed transformer-cnn architectures. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14388–14397, 2023. 2, 3

[34] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and Qianru Sun. Mnemonics training: Multi-class incremental learning without forgetting. *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 12245–12254, 2020. 2

[35] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11966–11976, 2022. 5, 2

[36] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017. 2

[37] Ming Lu and Zhan Ma. High-efficiency lossy image coding through adaptive neighborhood information aggregation. *arXiv preprint arXiv:2204.11448*, 2022. 2

[38] Ming Lu, Peiyao Guo, Huiqing Shi, Chuntong Cao, and Zhan Ma. Transformer-based image compression. *Data Compression Conference*, pages 469–469, 2022.

[39] 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*, 44(3):1247–1263, 2022. 2, 8

[40] D. Minnen and S. Singh. Channel-wise autoregressive entropy models for learned image compression. *Proceedings of the IEEE International Conference on Image Processing*, pages 3339–3343, 2020. 2

[41] D. Minnen, J. Ballé, and G. Toderici. Joint autoregressive and hierarchical priors for learned image compression. *Advances in Neural Information Processing Systems*, 31:10794–10803, 2018. 1, 2, 3, 5, 6, 4

[42] Guanbo Pan, Guo Lu, Zhihao Hu, and Dong Xu. Content adaptive latents and decoder for neural image compression. *Proceedings of the European Conference on Computer Vision*, pages 556–573, 2022. 2

[43] Yichen Qian, Zhiyu Tan, Xiuyu Sun, Ming Lin, Dongyang Li, Zhenhong Sun, Li Hao, and Rong Jin. Learning accurate entropy model with global reference for image compression. *International Conference on Learning Representations*, 2021. 2

[44] Yichen Qian, Xiuyu Sun, Ming Lin, Zhiyu Tan, and Rong Jin. Entroformer: A transformer-based entropy model for learned image compression. *International Conference on Learning Representations*, 2022. 2

[45] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5533–5542, 2017. 2, 4

[46] J. Rissanen and G. G. Langdon. Arithmetic coding. *IBM Journal of Research and Development*, 23(2):149–162, 1979. 3

[47] Sheng Shen, Huanjing Yue, and Jingyu Yang. Dec-adapter: Exploring efficient decoder-side adapter for bridging screen content and natural image compression. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12887–12896, 2023. 2

[48] Myungseo Song, Jinyoung Choi, and Bohyung Han. Variable-rate deep image compression through spatially-adaptive feature transform. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2360–2369, 2021. 2

[49] Kuan Tian, Yonghang Guan, Jinxi Xiang, Jun Zhang, Xiao Han, and Wei Yang. Effortless cross-platform video codec: A codebook-based method. *arXiv preprint arXiv:2310.10292*, 2023. 3

[50] Koki Tsubota, Hiroaki Akutsu, and Kiyoharu Aizawa. Universal deep image compression via content-adaptive optimization with adapters. *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2528–2537, 2023. 2

[51] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *arXiv preprint arXiv:2302.00487*, 2023. 1

[52] Yueqi Xie, Ka Leong Cheng, and Qifeng Chen. Enhanced invertible encoding for learned image compression. *Proceedings of the ACM International Conference on Multimedia*, pages 162–170, 2021. 2

[53] Fei Yang, Luis Herranz, Yongmei Cheng, and Mikhail G. Mozerov. Slimmable compressive autoencoders for practical neural image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4996–5005, 2021. 2

[54] Yibo Yang and Stephan Mandt. Computationally-efficient neural image compression with shallow decoders. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 530–540, 2023. 2

[55] Yibo Yang, Robert Bamler, and Stephan Mandt. Variational Bayesian quantization. *Proceedings of the International Conference on Machine Learning*, 119:10670–10680, 2020. 2

[56] Yibo Yang, Robert Bamler, and Stephan Mandt. Improving inference for neural image compression. *Advances in Neural Information Processing Systems*, 33:573–584, 2020. 2, 7

[57] Xi Zhang and Xiaolin Wu. Lvqac: Lattice vector quantization coupled with spatially adaptive companding for efficient learned image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10239–10248, 2023. 2

[58] Xiaosu Zhu, Jingkuan Song, Lianli Gao, Feng Zheng, and Heng Tao Shen. Unified multivariate gaussian mixture for efficient neural image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17591–17600, 2022. 4, 5

[59] Renjie Zou, Chunfeng Song, and Zhaoxiang Zhang. The devil is in the details: Window-based attention for image compression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17471–17480, 2022. 2, 3

# Towards Backward-Compatible Continual Learning of Image Compression

## Supplementary Material

## 7. Appendix: Model Architecture Details

In the main paper, Sec. 4.2 provides a high-level overview of the proposed model architecture. This section provides more details about the model architecture, such as the number of channels and stride sizes for each layer.

The detailed model architecture is shown in Fig. 7, where the model components are marked in the same way as in Sec. 4.2. Our model contains four phases, all of which have the same structure, while only different in (1) the number of feature channels, (2) the number of ConvNeXt blocks, and (3) the first phase starts from bias $e_0$ and $r_0$ instead of the feature maps from the previous phase.

The spatial dimensions (height and width) in the figure are for an input image with $256 \times 256$ pixels. Since the model is fully convolutional, the spatial dimensions of intermediate layer outputs scales accordingly with the input image size. Both initial bias features $e_0$ and $r_0$ have a shape of $1 \times 1 \times 128$, and they are repeated spatially to match the spatial dimensions of $z_1$.

## 8. Appendix: Training and Fine-tuning Details

Tab. 10 lists the pre-training and fine-tuning hyperparameters used in our experiments. For a fair comparison, we use the same hyperparameters for training all models, including our proposed model and the baseline models ( i.e., *MSH-VR* and *GMA-VR*). Note that the fine-tuning dataset varies for different sets of experiments. For data-incremental learning, we use CelebA-HQ [30], and for rate-incremental learning, we use COCO [32], which is the same as the pre-training dataset.

| | Pre-training | Fine-tuning |
|---|---|---|
| Data augmentation | Crop, h-flip | Crop, h-flip |
| Input size | 256x256 | 256x256 |
| Optimizer | Adam | Adam |
| Learning rate | $2 \times 10^{-4}$ | $1 \times 10^{-4}$ |
| LR schedule | Constant + cosine | Cosine |
| Weight decay | 0.0 | 0.0 |
| Batch size | 32 | 32 |
| # iterations | 500K | 100K |
| # images seen | 16M | 3.2M |
| Gradient clip | 2.0 | 2.0 |
| EMA | 0.9999 | - |
| GPU | $1 \times$ RTX 3090 | $1 \times$ A40 |
| Time | $\approx 51$ hours | $\approx 11$ hours |

Table 10. Training Hyperparameters. The GPU time is for training our proposed model, and all other hyperparameters are the same for all models.

## 9. Appendix: Variable-Rate Baseline Models

In the main paper (Sec. 5.2), we mentioned that we construct variable-rate versions of the two baseline models (*i.e.*, *MSH-VR* and *GMA-VR*) in order to use them in the rate-incremental learning experiment. Fig. 11 shows the rate-distortion performance of the variable-rate versions compared to the original ones. As shown in the figure, the variable-rate versions achieve similar performance as the original ones, which validates the our experimental setting.

## 10. Appendix: Experimental Results

### 10.1. PSNR-Bpp curves for the main experiments

Due to the space constraint, we show only BD-rate results without PSNR-bpp curves in the main paper. This section provides the PSNR-bpp curves for the main experiments (Sec. 5.3).

Fig. 8 shows the PSNR-bpp curves for data-incremental learning experiments, which includes the backward compatibility experiment (Fig. 8a) and the new-data performance experiment (Fig. 8b). For backward compatibility, it is clear that models with fine-tuned encoder and decoder suffer a significant performance drop on the old bitstreams, while other fine-tuned models obtain comparable performance as the pre-trained models. Among them, our proposed knowledge replay strategy achieves even better performance than using the pre-trained model directly. For new-data performance, our method achieves comparable performance as the models with fine-tuned encoder and decoder (which are not backward compatible), and outperforms the pre-trained models by a clear margin. These observations are consistent with what we have observed in the BD-rate results in the main paper.

We show the PSNR-bpp curves for rate-incremental learning experiments, including the *low-to-high* experiment (Fig. 9) and the *high-to-low* experiment (Fig. 10). The results are consistent with previous observations: (1) Fine-tuning the encoder and decoder does not preserve backward compatibility, while our approach does; and (2) Our approach even outperforms all other methods in terms of new-rate performance.

### 10.2. Fine-tuning the encoder does not generalize the model to new rates

We mentioned in Sec. 5.3 that fine-tuning the encoder alone cannot effectively extend the rate range of the pre-trained models. We provide an example for showing this in Fig. 12, where we show the rate-incremental learning (low → high)
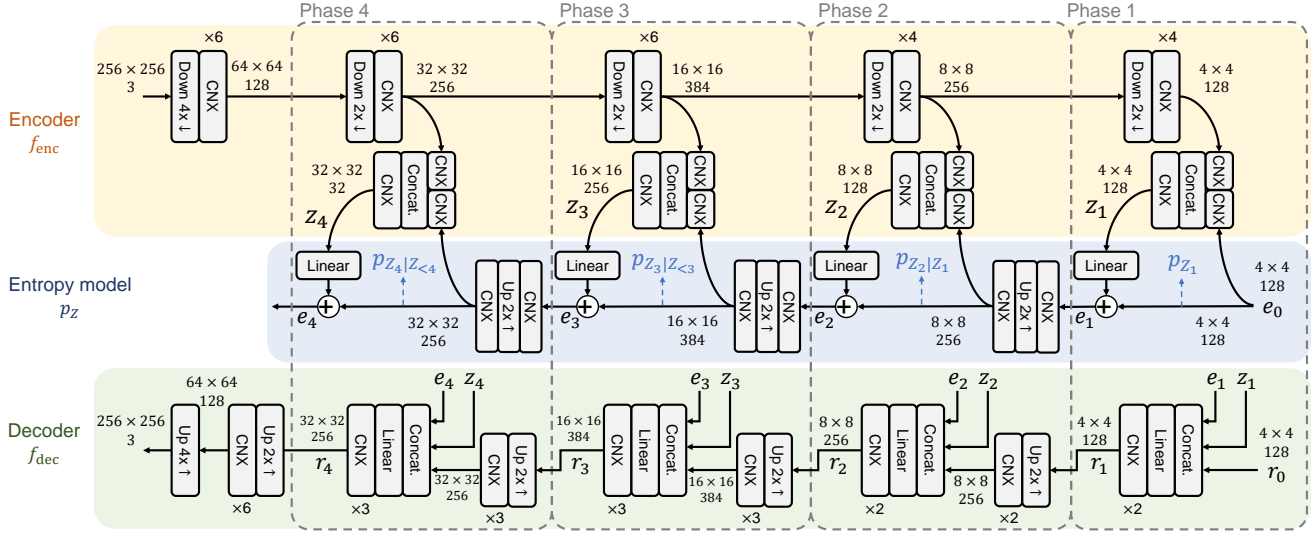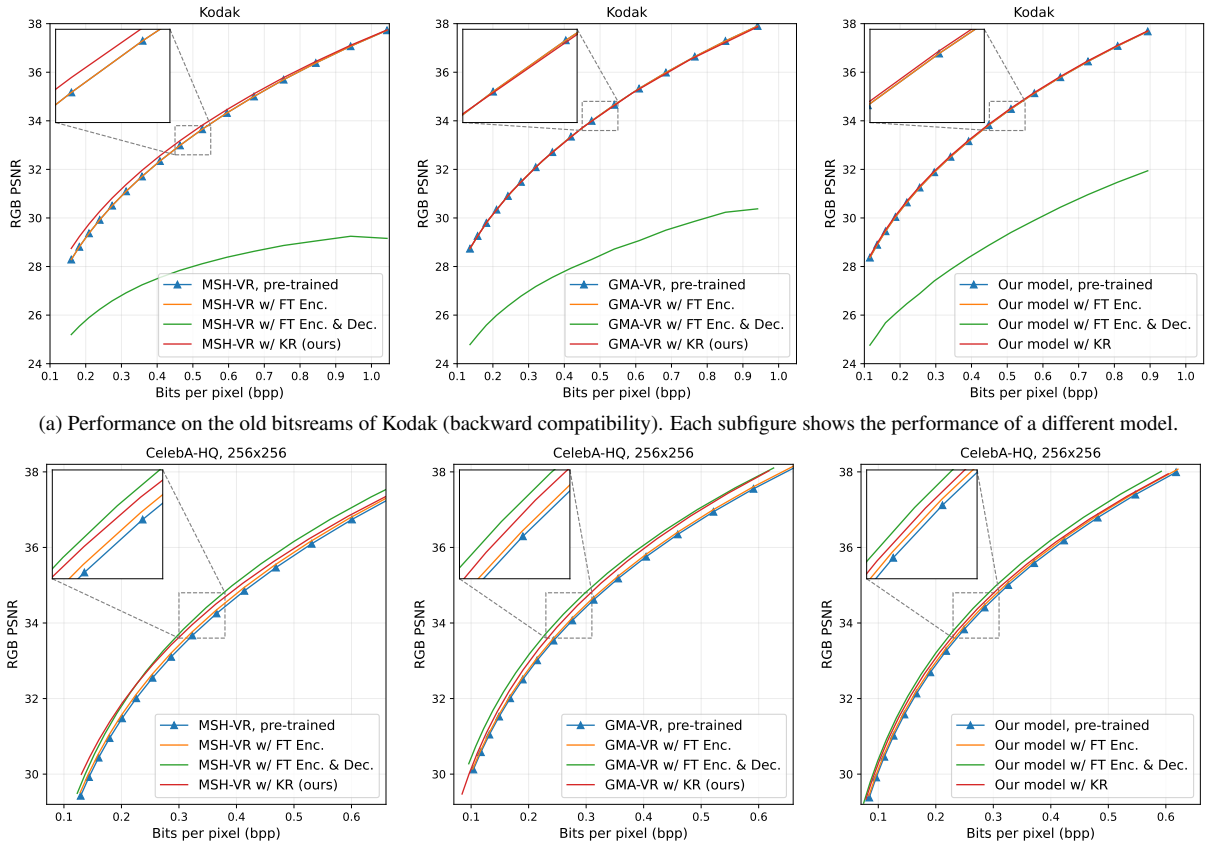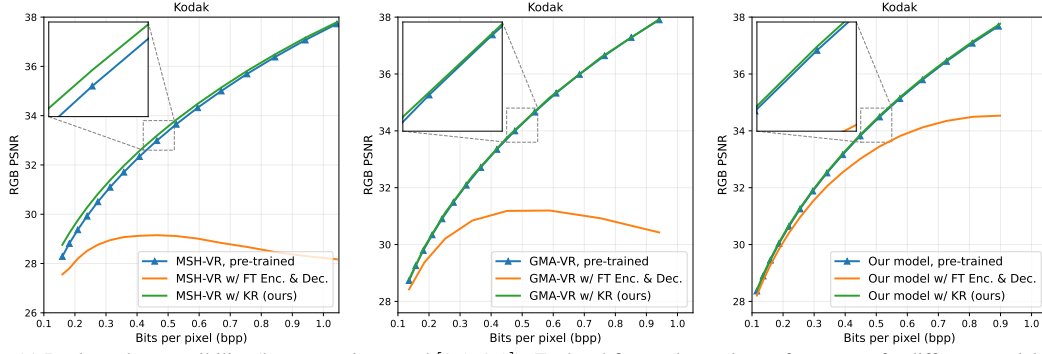
Figure 7. Detailed architecture of the proposed model. In the figure, *CNX* denotes a ConvNeXt block [35] conditioned on lagrange multiplier $\lambda$, as described in Fig. 6. Dimensionality of the layer outputs are shown in the format of *height × width* and *channels*, where the spatial dimensions (height and width) are for a $256 \times 256$ input image, and they scales linearly with the input image size.
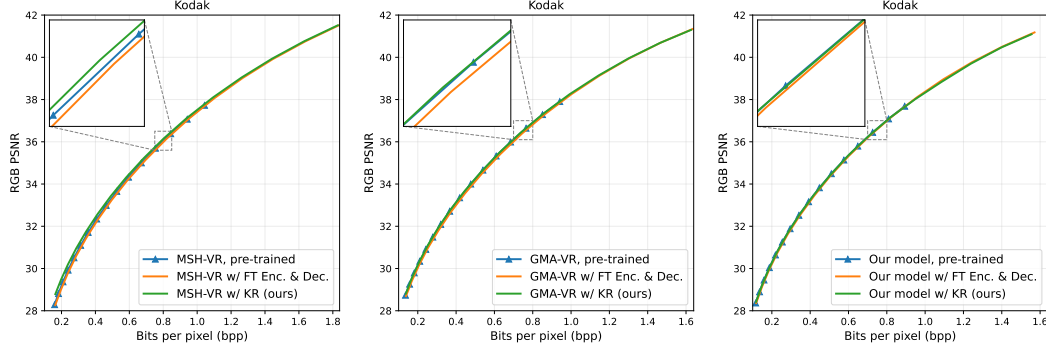


(a) Performance on the old bitsreams of Kodak (backward compatibility). Each subfigure shows the performance of a different model.



(b) Performance on CelebA-HQ (new-data performance). Each subfigure shows the performance of a different model.

Figure 8. PSNR-Bpp curves for **data-incremental learning** experiments. In figure (a), the *"models, pre-trained"* curves overlap with the *"models w/ FT Enc."* curves because their decoder are the same.
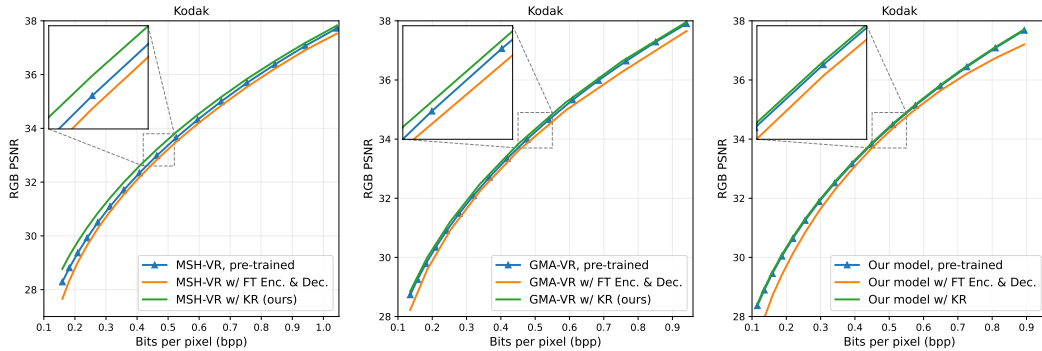
(a) Backward compatibility (bpp range is around $[0.1, 0.9]$). Each subfigure shows the performance of a different model.
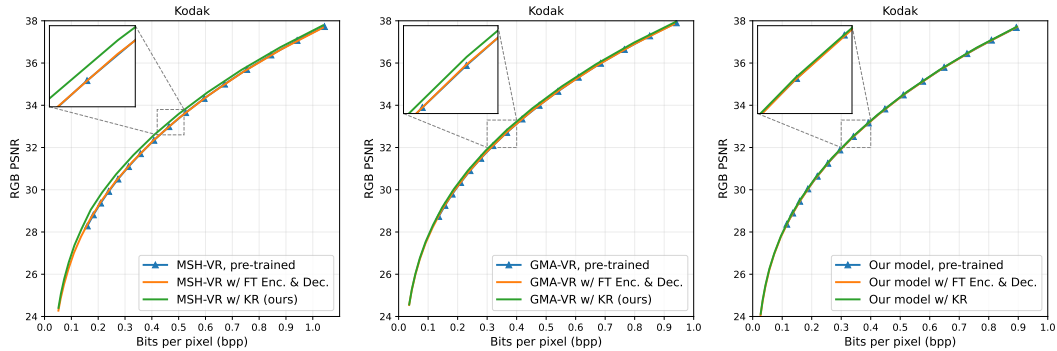


(b) New-rate performance (bpp range is around $[0.1, 1.6]$). Each subfigure shows the performance of a different model.

Figure 9. PSNR-Bpp curves for **rate-incremental learning (low → high)** experiments.



(a) Backward compatibility (bpp range is around $[0.1, 0.9]$). Each subfigure shows the performance of a different model.



(b) New-rate performance (bpp range is around $[0.03, 0.9]$). Each subfigure shows the performance of a different model.

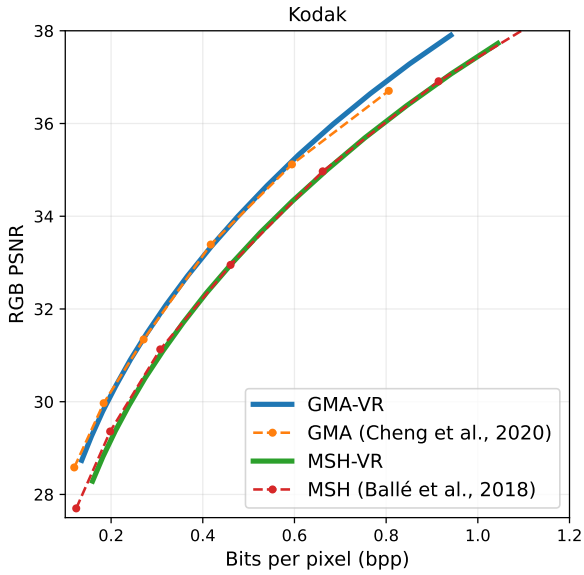Figure 10. PSNR-Bpp curves for **rate-incremental learning (high → low)** experiments.

Figure 11. The variable-rate version of the baseline models that we constructed (MSH-VR and GMA-VR) are comparable to the original ones (MSH [41] and GMA [11]) in terms of PNSR-bpp performance on Kodak.
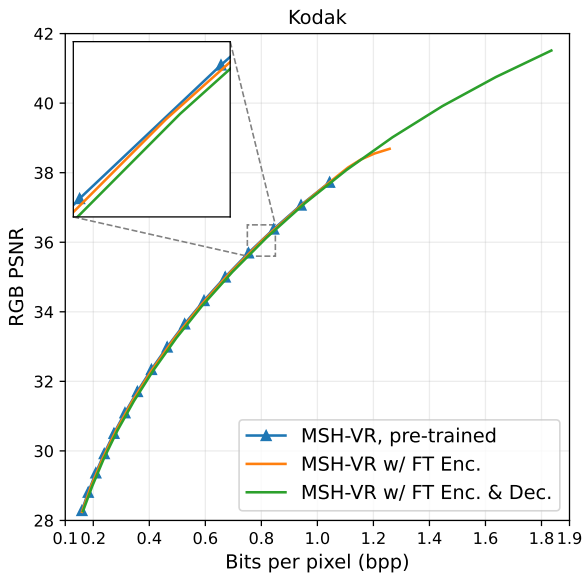


Figure 12. Fine-tuning the encoder does not effectively generalize the pre-train model (MSG-VR, for example) to new rates.

performance of the pre-trained MSH-VR, the one with fine-tuned encoder (*MSH-VR w/ FT Enc.*), and the one with fine-tuned encoder and decoder (*MSH-VR w/ FT Enc. & Dec.*). As shown in the figure, fine-tuning the encoder marginally extends the rate range of the pre-trained model, and the PSNR drops visibly when the rate is higher than maximum rate of the pre-trained model. Thus, we do not use this strategy in our rate-incremental learning experiments.