

THEORETICAL BOUND-GUIDED HIERARCHICAL VAE FOR NEURAL IMAGE CODECS

Yichi Zhang, Zhihao Duan, Yuning Huang, Fengqing Zhu

Elmore Family School of Electrical and Computer Engineering,
Purdue University, West Lafayette, Indiana, U.S.A.
{zhan5096, duan90, huan1781, zhu0}@purdue.edu

ABSTRACT

Recent studies reveal a significant theoretical link between variational autoencoders (VAEs) and rate-distortion theory, notably in utilizing VAEs to estimate the theoretical upper bound of the information rate-distortion function of images. Such estimated theoretical bounds substantially exceed the performance of existing neural image codecs (NICs). To narrow this gap, we propose a theoretical bound-guided hierarchical VAE (BG-VAE) for NIC. The proposed BG-VAE leverages the theoretical bound to guide the NIC model towards enhanced performance. We implement the BG-VAE using Hierarchical VAEs and demonstrate its effectiveness through extensive experiments. Along with advanced neural network blocks, we provide a versatile, variable-rate NIC that outperforms existing methods when considering both rate-distortion performance and computational complexity. The code is available at [BG-VAE](#).

Index Terms— Lossy Image Compression, Knowledge Distillation, Hierarchical VAE

1. INTRODUCTION

Lossy image compression, an important problem in image processing, aims to compress images to a low-rate representation while retaining high reconstruction quality. It is essential for efficient storage and transmission in numerous applications, from web media to satellite imagery. With the rapid advances in deep learning, the landscape of lossy image compression has undergone a significant transformation. Recent works not only achieved substantial progress in practical compression methods but also deepened the theoretical analysis. The majority of research in this field focuses on the practical aspect, where various neural network architectures [1, 2] and probability models [3, 4] have been proposed to improve the rate-distortion (R-D) performance. This leads to a series of strong neural image codecs (NICs) [2, 5, 6] that outperform traditional codecs such as H.266/VVC [7].

Despite the impressive R-D performance of recent NICs, it has been shown that there is still a considerable gap between their performance and theoretically achievable limits [8, 9]. Recall that the best achievable R-D performance for a data source is described by its *information rate-distortion (R-D) function*, a fundamental quantity in lossy compression. By

bounding the information R-D function using deep variational autoencoders [10] (VAEs), Yang and Mandt [8] show that NICs could be improved by at least +1 dB in PSNR at various rates. Interestingly, the neural network models used for bounding the information R-D function closely resemble recent NICs in the sense that both are hierarchical VAEs.

Inspired by this resemblance, this paper explores the possibility of leveraging the estimated bound information R-D functions to improve the NICs. Specifically, we leverage the teacher-student framework [11] that is prevalent in knowledge distillation methods. The NIC’s theoretical bound naturally serves as a *teacher*, while the NIC itself takes the role of the *student*. During training, the NIC is guided by a theoretical bound model, thereby enhancing its efficacy. The similarity in model structure and model size between the bound and the NIC obviates the need for a typically larger teacher model in the teacher-student framework. The availability of theoretical bound simplifies the teacher selection process and reduces training resource consumption. Furthermore, we exemplify a hierarchical NIC named BG-VAE (theoretical Bound-Guided VAE) with well-designed modules. Extensive experiments demonstrate the superior performance of our method compared to existing approaches.

Our contributions are summarized as follows:

- We propose a teacher-student framework that uses theoretical bounds to guide the training of NICs.
- We present new network modules and construct an efficient hierarchical model that leverages spatial and spectral information.
- Overall, we present BG-VAE, a NIC framework including the teacher-student training strategy and model architectures. Extensive experiments are conducted to demonstrate the effectiveness of the proposed method.

2. RELATED WORK

In this section, we briefly review previous works and summarize the preliminaries.

2.1. Neural Image Codecs and VAEs

Most existing NICs follow the scheme of *transform coding*, where images are transformed to a latent space for decorrelation and energy compression, followed by quantization

and entropy coding. Early works stack convolutions to parameterize the codec [12, 13]. Later, attention mechanisms [3] and transformers [14, 15], are developed to improve the performance. Another line of research focuses on developing context models to aid entropy coding involved with leveraging hierarchical latent variables [13, 16, 4] and exploring correlations between pixels [17] as well as between channels of latent variables [2]. Essentially, the framework of these methods is similar to a one-layer [12] or two-layer [13] VAEs.

VAEs [10] is a critical class of latent variable models, especially effective for complex data like images. In VAEs, data X and latent variable Z are linked through the joint distribution $p_{X,Z}(x, z) = p_{X|Z}(x|z) \cdot p_Z(z)$. VAEs feature an encoder-like approximate posterior $q_{Z|X}$ and a decoder-like $p_{X|Z}$, framing them as stochastic autoencoders.

In the NIC scenarios, VAEs employ a deterministic decoder f_{dec} , leading to a lossy reconstruction $\hat{X} = f_{\text{dec}}(Z)$. For high-dimensional data, hierarchical VAEs [18] enhance flexibility and expressiveness. They employ a series of latent variables $Z_{1:N} \triangleq \{Z_1, \dots, Z_N\}$ in an autoregressive manner: $p_{Z_{1:N}} = \prod_{i=1}^N p_{Z_i|Z_{<i}}$. The architecture progresses from low to high dimensions, capturing the image’s granular details. The hierarchical VAE-based NIC training loss extends the standard single-layer loss across multiple latent variables:

$$\mathcal{L}_\lambda = \mathbb{E}_{X \sim p_{\text{data}}, Z_{1:N} \sim q_{1:N}} \left[\sum_{i=1}^N D_{\text{KL}}(q_i || p_i) + \lambda \cdot d(X, \hat{X}) \right], \quad (1)$$

q_i and p_i denote the posterior and prior of each latent variable, approximated via ancestral sampling.

As for its theoretical bounds, Yang *et al.* [8] established an upper bound on the information R-D function for image sources by training VAEs. In doing so, they show that there still exists a sizeable room for improvement over current NICs (+1 dB in PSNR at various rates). Duan *et al.* [9] reported an improved upper bound on the information R-D function of images by extending it to hierarchical VAE and variable rate image compression, which showed that at least 30% BD-rate reduction w.r.t. the VVC codec is achievable.

2.2. Knowledge Distillation

Knowledge Distillation (KD), initially introduced by Hinton *et al.* [19], balances model performance and efficiency. It serves as a training strategy, enhancing the performance of lightweight networks through knowledge transfer from a high-capacity teacher model, predominantly used in image classification [19] and object detection [11]. In contrast, KD in NIC remains relatively unexplored. Fu *et al.* [20] developed an improved three-step KD training scheme for balancing decoder network complexity and performance, transferring both final and intermediate outputs from the teacher to the student network. However, this approach still relies on a large teacher model and a less effective MSE loss function.

Remarks. Prior efforts in NIC predominantly emphasize model designs to achieve commendable performance. However, another effective way may involve identifying the theoretical bound of NIC and utilizing the bound to improve practical performance. Our BG-VAE begins with bound-guided training, forcing the practical NIC (*student*) to mimic the bound (*teacher*) behavior. Due to the similarity in model structure and model size of the bound and NIC, this method bypasses the need for a larger teacher model, as is customary in the KD method, thereby reducing training complexity.

3. PROPOSED METHOD

An overview of the proposed BG-VAE is presented in Fig. 1. We begin by describing the bound-guided training and the loss function, as elaborated in Section 3.1. Subsequently, the neural network architecture employed for implementing our BG-VAE is described in Section 3.2.

3.1. Theoretical Bound-Guided Training

As depicted in Fig. 1(a), our training recipe involves two models: a model that is used for estimating theoretical bounds, denoted by B , and the NIC model for practical deployment, denoted by P . We refer to the former model as the *teacher* and the latter as *student*. The teacher model’s implementation adheres to the approaches from [9]. Both models share a similar network structure as shown in Fig 1(b), and P is trained to mimic B ’s behavior including the intermediate features and the reconstructed image.

3.1.1. Affinity Matrix-based Feature Alignment

Precise similarity measurement is crucial in feature alignment, especially given the unbounded feature representation space in regression problems [21], a similar case for NICs. This contrasts with the more constrained feature spaces in classification problems. As a result, distillation methods effective in classification [19] are not suitable for compression. Inspired by FAKD [21], we develop an affinity matrix-based feature alignment process.

We extract feature maps at two points: after downsampling (*Wavelet Down* in Fig. 1(b)) and before the upsampling (*Wavelet Up* in Fig. 1(b)). For the teacher model, the feature map $F^{(B)} \in \mathbb{R}^{b \times c \times \frac{w}{n} \times \frac{h}{n}}$ is extracted directly, with b, c, w, h, n , represent the batch size, channels, width, height, the downsampled ratio at that stage, respectively. For the student model, we first use two B-ConvNeXt blocks (refer to Sections 3.2) to transform the feature map, and subsequently, we get $F^{(P)} \in \mathbb{R}^{b \times c \times \frac{w}{n} \times \frac{h}{n}}$. Notably, we employ an ensemble strategy [22] for the two B-ConvNeXt blocks.

Next, we compute the affinity matrix $A^{(B)}, A^{(P)}$ using $F^{(B)}$ and $F^{(P)}$. Our method emphasizes the spatial relationships within $A^{\{(B),(P)\}} \triangleq A$. This is achieved by normalizing the matrix along the channel dimension and excluding this dimension through matrix multiplication during the computation. The affinity matrix A is thus formulated by considering the spatial interrelations between pixels, highlight-

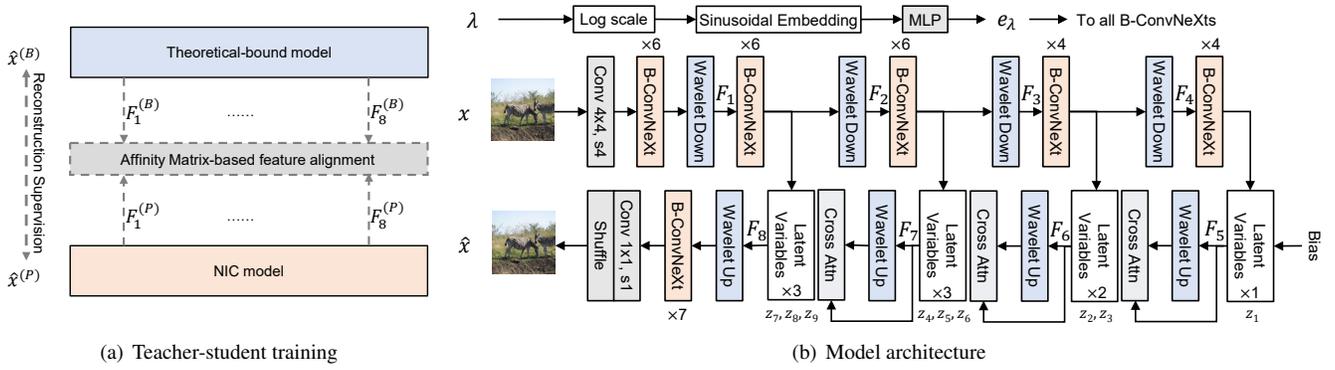


Fig. 1. An overview of the proposed BG-VAE. (a) the bound-guided framework, (b) the model used to implement BG-VAE.

ing their mutual dependencies and interactions. The affinity matrix $A \in \mathbb{R}^{b \times \frac{hw}{n^2} \times \frac{hw}{n^2}}$ is determined by:

$$A = \left(\frac{F}{\sqrt{\sum_{k=1}^c F_k^2 + \varepsilon}} \right)^T \times \frac{F}{\sqrt{\sum_{k=1}^c F_k^2 + \varepsilon}} \quad (2)$$

where F is reshaped to $b \times c \times \frac{hw}{n^2}$ first. F_k is the k -th channel of F . ε is $1e-8$ used to improve numerical stability.

To assess the similarity between $A^{(B)}$ and $A^{(P)}$, we utilize two indicators, the L1 loss for element-by-element difference assessment, and the cosine similarity loss, which evaluates the angular relationship between matrices, offering insights beyond mere element-wise comparison. Consequently, the feature affinity-based guidance loss $\mathcal{L}_{\text{feature}}$ is expressed as:

$$\mathcal{L}_{\text{feature}} = \underbrace{\frac{1}{n} \sum_{i=1}^n |A_i^{(B)} - A_i^{(P)}|}_{\text{L1 Loss}} + 1 - \underbrace{\frac{A^{(B)} \cdot A^{(P)}}{\|A^{(B)}\| \|A^{(P)}\|}}_{\text{Cosine Similarity Loss}} \quad (3)$$

where $A_i^{(B)}$ and $A_i^{(P)}$ are the i -th element of the teacher and student model's affinity matrix respectively.

3.1.2. Overall Loss

Our approach includes the teacher model's supervisory role over the output of the practical codec. This is achieved by calculating the reconstruction supervision loss, denoted as \mathcal{L}_{rs} :

$$\mathcal{L}_{\text{rs}} = \frac{1}{n} \sum_{i=1}^n \left| \hat{x}_i^{(B)} - \hat{x}_i^{(P)} \right| \quad (4)$$

$\hat{x}_i^{(B)}$ and $\hat{x}_i^{(P)}$ are the i -th elements of the reconstructed image from the teacher model and the student model, respectively. Consequently, the final loss function is as follows:

$$\mathcal{L} = \mathcal{L}_{\lambda} + w_1 \cdot \sum_{j=1}^8 \mathcal{L}_{\text{feature}_j} + w_2 \cdot \mathcal{L}_{\text{rs}} \quad (5)$$

where w_1 and w_2 are weighting factors adjusting magnitude orders which are set to 1 in our method. $\mathcal{L}_{\text{feature}_j}$ represents the

loss calculated by using F_j . Variable-rate training is achieved by \mathcal{L}_{λ} (Eq. (1)), where λ is randomly sampled from the range [64, 8192] during training (this range approximately corresponds to [0.2, 2.2] bits per pixel (bpp)). The entire model is conditioned on λ using the λ -embedding network (shown at the top of Fig. 1(b)) and through B-ConvNeXt.

At test time, the coding rate can be adjusted by tuning the input λ , effectively achieving variable-rate compression.

3.2. Model Architectures

This section presents a neural network implementation of the BG-VAE framework and its components in detail.

As shown in Fig 1(b), to build the BG-VAE, we utilized a hierarchical VAE-based structure transmitting 9 latent variables (z_1 to z_9) in four stages. In the top-to-bottom path, ratio 4 downsampling is first achieved using a 4×4 , stride 4 convolution, followed by feature extraction using B-ConvNeXt and Wavelet Downsampling for ratio 2 downsampling. Latent Variable Blocks transmit latent variables z_i . In the bottom-to-top path, Wavelet Upsampling is employed for ratio 2 upsampling. To improve global information integration, we integrate a cross-attention module [23], which harnesses low-resolution global information, thereby facilitating a more comprehensive global view. Finally, 1×1 , stride 1 convolution, and shuffle enable ratio 4 upsampling. Each module will be elaborated in the following.

Balanced ConvNeXt block. Based on ConvNeXt block [24, 25], we introduce the Balanced ConvNeXt block (B-ConvNeXt), depicted in Fig. 2. While retaining the core architecture of ConvNeXt, we implement key modifications: (1): After depth-wise convolution, the DC (direct current component) is extracted by averaging each channel and subtracted from the original features to isolate the HC (high-frequency component). Two learnable parameters α, β are then applied to modulate the balance between the DC and HC before re-integrating them with the original features; (2): We enable variable rates by conditioning on e_{λ} , an embedded output from the λ embedding network (illustrated at the top of Fig. 1(b)). e_{λ} passes through the GELU function and a Linear layer to scale features after LayerNorm in the block, making them conditional on λ .

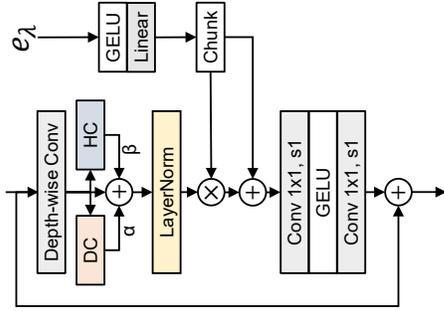


Fig. 2. The structure of Balanced ConvNeXt block.

Wavelet Up/Down Sampling. We utilize the Discrete Wavelet Transform (DWT) for up/down sampling. DWT divides data into four subbands: the ILL, capturing coarse object structures, and the ILH, IHL, and IHH, detailing fine textures. These subbands are merged using a single convolution, reducing dimensions. Before applying the inverse DWT (IDWT), another single convolution is used to re-establish the original dimensions. DWT and IDWT are invertible operations, which preserves feature quality and fidelity.

Latent Variable Blocks. The Latent Variable Block is shown in Fig. 3. The left-hand side of Fig. 3 shows the posterior branch of the latent variable, and the middle part indicates the prior branch. The posterior branch consists of three B-ConvNeXt blocks, a concatenation operation, and two convolutions. The prior branch, in contrast, consists of a single convolution that facilitates decoding. The structure of the Latent Variable Block for the teacher model is the same, except for the distribution of Posterior and Prior is Gaussian.

Posteriors. The posterior distribution of Z_i given x and $z < i$ is defined as:

$$q_i \triangleq U\left(\mu_i - \frac{1}{2}, \mu_i + \frac{1}{2}\right) \quad (6)$$

where μ_i is the output of the posterior branch, depending on the image x and preceding latent variables $z_{<i}$. Once q_i is obtained, z_i is sampled as $z_i \leftarrow \mu_i + u$, where u is a random sample from $U(-\frac{1}{2}, \frac{1}{2})$ during training, and during testing, it is replaced with scalar quantization.

Priors. The prior distribution p_i is defined as a conditional Gaussian convolved with a uniform distribution:

$$p_i \triangleq \mathcal{N}(\hat{\mu}_i, \hat{\sigma}_i^2) * U\left(-\frac{1}{2}, \frac{1}{2}\right) \quad (7)$$

where $\mathcal{N}(\hat{\mu}_i, \hat{\sigma}_i^2)$ represents the Gaussian probability density function. The mean $\hat{\mu}_i$ and standard deviation $\hat{\sigma}_i$ are predicted by the prior branch. The probability mass function (PMF) P_i is then defined as:

$$P_i(n) \triangleq p_i(\hat{\mu}_i + n | z_{<i}), n \in \mathbb{Z}. \quad (8)$$

which is used for the entropy coding/decoding of z_i .

For details on the model architecture, please refer to the Supplementary Information.

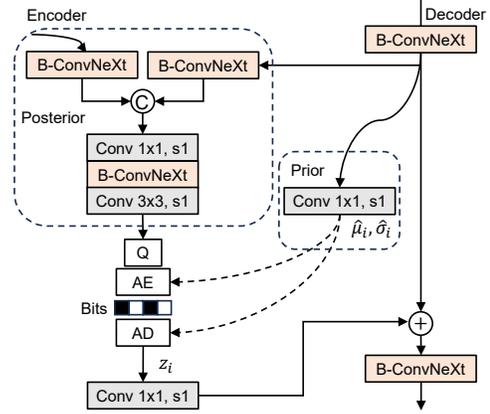


Fig. 3. Illustration of the i -th Latent Variable Block.

4. EXPERIMENTS

4.1. Experimental Settings

Training. We use COCO2017 [26] as our training sets. It contains 118,287 images with a resolution of 640×420 pixels. We randomly cropped 256×256 patches for training. The model is trained for 2M iterations with a batch size of 32 and a learning rate of $2e^{-4}$. For all ablation experiments, we train the models for 500k iterations. More details are shown in the supplementary information. B is pre-trained using the same settings and fixed during our BG-VAE training.

Testing. Three widely used benchmark datasets, including Kodak, Tecnick, and CLIC 2022, are used to evaluate the performance of the proposed method.

4.2. Quantitative Results

We compare our proposed method with prevalent NICs including fixed rate methods: M&S Hyperprior [13], Cheng2020 [3], STF [1], ELIC [2], TCM-S [5]; variable rate method: QARV [4]; and rule-based method: VVC [7]. We use VTM-18.0 All Intra as the anchor to calculate BD-Rate.

Table 1 reports the BD-Rate reduction of each method against the VVC anchor on three datasets. Our BG-VAE achieves impressive performance on each dataset, -7.04% BD-Rate on Kodak, -8.21% BD-Rate on Tecnick, and -6.33% BD-Rate on CLIC 2022. Fig. 4 further plots RD curves of all methods. Our BG-VAE performs best on Kodak and Tecnick datasets and gains lower than TCM-S on CLIC2022. However, the total parameters of our BG-VAE are only 36% of TCM-S, highlighting our high efficiency especially when considering our BG-VAE covers a wide range of bpp (0.2 bpp to 2.2 bpp).

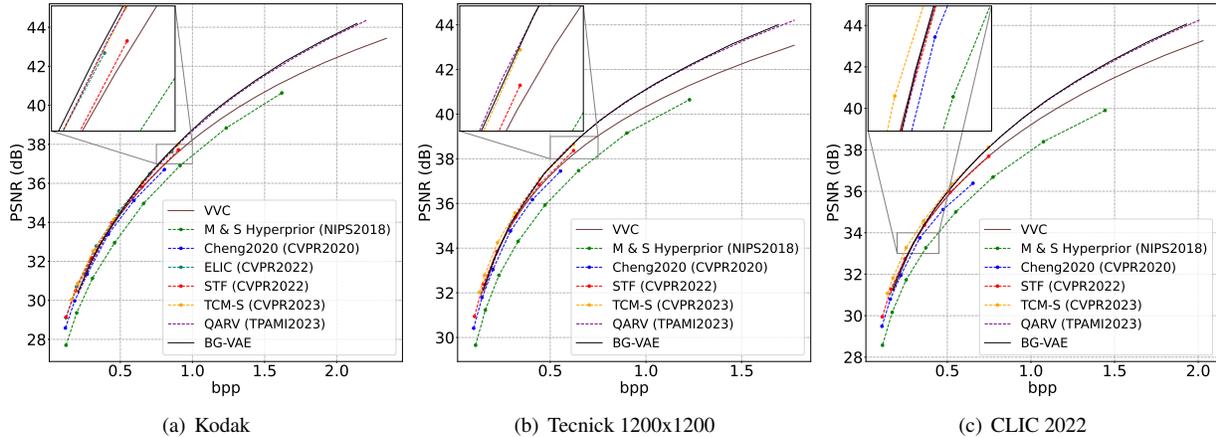
4.3. Complexity

We measure the computational complexity using the total parameters (The total parameters of the fixed rate method are obtained by summing up model parameters for all bpp.), encoding time (Enc.), and decoding time (Dec.) on CPU and GPU, as shown in Table 1. In comparison to other methods, our model's total parameter is notably efficient, being

Table 1. Computational Complexity and BD-Rate Compared to Existing Learning-based Methods

Method	Total Params.	Latency (CPU)		Latency (GPU)		BD-Rate (%) w.r.t. VTM 18.0		
		Enc.	Dec.	Enc.	Dec.	Kodak	Tecnick	CLIC2022
M & S Hyperprior [17]	98.4M	0.759s	0.830s	0.033s	0.030s	17.06%	26.09%	29.11%
Cheng2020 [3]	115.3M	3.605s	6.081s	1.048s	2.376s	3.89%	5.92%	8.32%
STF [1]	599.1M	2.373s	2.738s	0.076s	0.068s	-2.55%	-2.35%	-1.17%
ELIC* [2]	202.8M	2.131s	2.187s	0.125s	0.062s	-5.69%	-	-
TCM-S [5]	271.1M	<u>0.837s</u>	0.945s	0.095s	0.086s	-5.54%	<u>-8.15%</u>	-8.52%
QARV [4]	93.4M	0.852s	0.309s	0.098s	0.068s	-5.82%	-7.79%	-6.13%
BG-VAE	<u>97.4M</u>	0.990s	<u>0.437s</u>	<u>0.082s</u>	<u>0.055s</u>	-7.04%	-8.21%	-6.33%

Test Conditions: Intel(R) Core(TM) i7-12700K CPU, Nvidia 3090 GPU. The enc./dec. time is averaged over all 24 images in Kodak, including entropy enc./dec. time. *: We reproduced ELIC [2] to calculate the runtime. ELIC results on Kodak are officially provided. **Bold** and underlined indicate the best and the second best, respectively.

**Fig. 4.** RD curves of various methods. Please zoom in for more details.

the second smallest and exceeding QARV by only 4M parameters. Despite this minimal difference, our model achieves a 1% higher BD-Rate Reduction. While M&S Hyperprior’s parameter is similar, its performance on the Kodak lags behind our results by 24% BD-Rate. In contrast, the other methods have a substantially higher number of parameters. Regarding latency, BG-VAE’s encoding time is comparable with other methods and the decoding time is only marginally longer than QARV on the CPU and longer than M&S Hyperprior on the GPU, which demonstrates the high efficiency of BG-VAE.

4.4. Ablation Study

A series of ablation studies are conducted to verify the contribution of the bound-guided method and each module.

Bound-Guided Training: One of the key components of our proposed method is the bound-guided training, as shown in Table 2. Eliminating the bound-guided training, see “BG-VAE base model”, results in an average increase of 0.99% BD-Rate compared to “BG-VAE”, highlighting the bound-guided training’s effectiveness. Further, replacing the bound model with a larger teacher model (LM), as shown in “W / LM”, demonstrates the efficiency of using the bound to guide the NIC. Notably, due to the larger number of parameters of LM (124.5M vs 97.4M), training the LM incurs greater resource expenditure than the bound. Additionally, maintaining

the bound model while replacing the $\mathcal{L}_{\text{feature}}$, \mathcal{L}_{rs} with \mathcal{L}_{MSE} (“W / \mathcal{L}_{MSE} ”) further underscores the affinity matrix’s superiority. In BG-VAE, features F_1 to F_8 , depicted in Fig 1(b) and Section 3.1.2, are used to calculate $\mathcal{L}_{\text{feature}}$. It’s also noteworthy that limited supervision, using only F_5 to F_8 (“W / F_5 to F_8 ”), proves less effective.”

Table 2. Ablation experiments on bound-guided framework

Settings	Kodak	Tecnick	CLIC2022
BG-VAE base model	-3.58%	-3.88%	-2.30%
W / LM	-3.79%	-4.12%	-3.09%
W / \mathcal{L}_{MSE}	-3.97%	-4.28%	-3.14%
W / F_5 to F_8	-3.92%	-4.52%	-3.12%
BG-VAE	-4.46%	-4.78%	-3.50%

Network Architecture: we conducted individual assessments of three advanced modules in Table 3. Removing Wavelet sampling (“W / o Wavelet sampling”) had negligible impact on Kodak images but significantly affected high-resolution Tecnick and CLIC2022 images, likely due to their richer textural details, suggesting its importance in handling textural details due to the invertibility of Wavelet downsampling, unlike convolution downsampling which loss details. Additionally, eliminating the Balancing factor (“W / o Bal-

ancing factor”) which entails using the ConvNeXt block in place of the B-ConvNeXt block, uniformly decreased performance across all datasets, underscoring its effectiveness with a simple design. Lastly, our experiment “W / o Cross-Attn” demonstrates its integral role across all three datasets. This confirms the effectiveness of Cross-Attention in leveraging global information to enhance overall performance.

Table 3. Ablation experiments on model architecture

Settings	Kodak	Tecnick	CLIC2022
W / o Wavelet sampling	-3.53%	-1.80%	-1.21%
W / o Balancing factor	-3.11%	-2.45%	-1.15%
W / o Cross-Attn	-3.03%	-3.09%	-1.79%
BG-VAE base model	-3.58%	-3.88%	-2.30%

5. CONCLUSION

In this paper, we present a theoretical Bound-Guided hierarchical VAE (BG-VAE) for NIC. The core of BG-VAE is its bound-guided training, which effectively boosts performance by utilizing the theoretical bound to guide the NIC model in training. Additionally, we develop several efficient modules that adeptly harness spatial-spectral information, forming the backbone of BG-VAE. These designs together present a versatile, variable-rate compression method. Extensive experimental results underscore BG-VAE’s superior performance when compared to the traditional rule-based VVC and other neural image codecs, marking a significant advancement in the field of image compression.

6. REFERENCES

- [1] 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*, pp. 17492–17501, June 2022.
- [2] 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*, pp. 5718–5727, June 2022.
- [3] 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*, pp. 7936–7945, June 2020.
- [4] 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*, vol. 46, no. 1, pp. 436–450, 2024.
- [5] 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*, pp. 14388–14397, June 2023.
- [6] Yichi Zhang, Zhihao Duan, Ming Lu, Dandan Ding, Fengqing Zhu, and Zhan Ma, “Another way to the top: Exploit contextual clustering in learned image coding,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, pp. 9377–9386, Mar. 2024.
- [7] Benjamin Bross, Ye-Kui Wang, Yan Ye, Shan Liu, Jianle Chen, Gary J Sullivan, and Jens-Rainer Ohm, “Overview of the versatile video coding (VVC) standard and its applications,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 10, pp. 3736–3764, 2021.
- [8] Yibo Yang and Stephan Mandt, “Towards empirical sandwich bounds on the rate-distortion function,” *International Conference on Learning Representations*, Apr. 2022.
- [9] Zhihao Duan, Jack Ma, Jiangpeng He, and Fengqing Zhu, “An improved upper bound on the rate-distortion function of images,” pp. 246–250, 2023.
- [10] D. Kingma and M. Welling, “Auto-encoding variational bayes,” *International Conference on Learning Representations*, Apr. 2014.
- [11] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao, “Knowledge Distillation: A Survey,” *International Journal of Computer Vision*, vol. 129, pp. 1789–1819, 2021.
- [12] Johannes Ballé, Valero Laparra, and Eero P. Simoncelli, “End-to-end optimized image compression,” *International Conference on Learning Representations*, Apr. 2017.
- [13] J. Ballé, D. Minnen, S. Singh, S. Hwang, and N. Johnston, “Variational image compression with a scale hyperprior,” *International Conference on Learning Representations*, Apr. 2018.
- [14] Ming Lu and Zhan Ma, “High-efficiency lossy image coding through adaptive neighborhood information aggregation,” *arXiv preprint arXiv:2204.11448*, Oct. 2022.
- [15] Yichi Zhang, Dandan Ding, Zhan Ma, and Zhu Li, “A reconfigurable framework for neural network-based video in-loop filtering,” *ACM Transactions on Multimedia Computing, Communications and Applications*, 2024.
- [16] Yueyu Hu, Wenhan Yang, and Jiaying Liu, “Coarse-to-Fine Hyper-Prior Modeling for Learned Image Compression,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, pp. 11013–11020, 2020.
- [17] D. Minnen, J. Ballé, and G. Toderici, “Joint autoregressive and hierarchical priors for learned image compression,” *Advances in Neural Information Processing Systems*, vol. 31, pp. 10794–10803, Dec. 2018.
- [18] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling, “Improved variational inference with inverse autoregressive flow,” *Advances in Neural Information Processing Systems*, vol. 29, Dec. 2016.
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, “Distilling the Knowledge in a Neural Network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [20] Haisheng Fu, Feng Liang, Jie Liang, Yongqiang Wang, Guohe Zhang, and Jingning Han, “Fast and High-Performance Learned Image Compression With Improved Checkerboard Context Model, Deformable Residual Module, and Knowledge Distillation,” *arXiv preprint arXiv:2309.02529*, 2023.
- [21] Zibin He, Tao Dai, Jian Lu, Yong Jiang, and Shu-Tao Xia, “Fakd: Feature-Affinity Based Knowledge Distillation for Efficient Image Super-Resolution,” *2020 IEEE International Conference on Image Processing*, pp. 518–522, 2020.
- [22] Yudong Chen, Sen Wang, Jiajun Liu, Xuwei Xu, Frank de Hoog, and Zi Huang, “Improved feature distillation via projector ensemble,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 12084–12095, 2022.
- [23] Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda, “CrossViT: Cross-Attention Multi-Scale Vision Transformer for Image Classification,” *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 357–366, 2021.
- [24] Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang, “Anti-Oversmoothing in Deep Vision Transformers via the Fourier Domain Analysis: From Theory to Practice,” *arXiv preprint arXiv:2203.05962*, 2022.
- [25] 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*, pp. 11976–11986, 2022.
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft COCO: Common Objects in Context,” *European Conference on Computer Vision*, pp. 740–755, 2014.

Supplementary Information

This document provides more details about our proposed method and comparison.

1.1. Detailed Architecture

1.1.1. Wavelet Up/Down Sampling

For downsampling, features are first decomposed into four subbands, and their resolution is reduced by the DWT, followed by a 1×1 convolution to reduce the channel dimension. For upsampling, on the other hand, they are first expanded back to their original channel dimensions and then we perform the IDWT, as shown in Fig. 1.

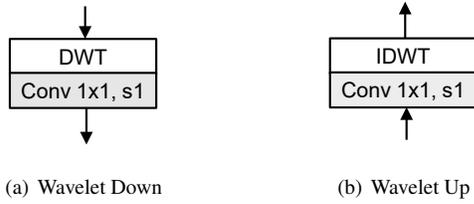


Fig. 1. The structure of Wavelet Up/Down sampling.

1.1.2. Cross-Attention

The cross-attention module [23] is illustrated in Fig. 2(a), and the attention mechanism is detailed in Fig. 2(b). This module [23] is designed to harness low-resolution global information, thereby facilitating a more comprehensive global view. F_H represents high-resolution features and F_L denotes low-resolution features from the previous stage. CPE (conditional positional encoding) is utilized [27].

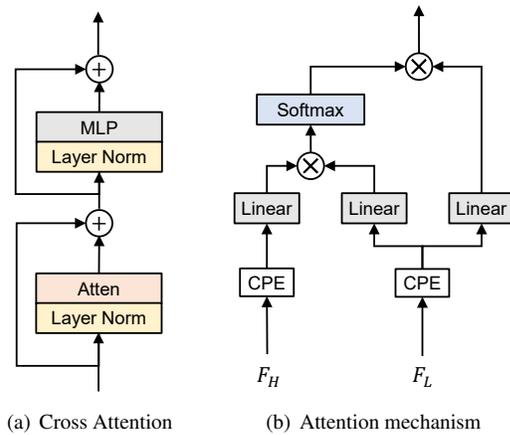


Fig. 2. The structure of Cross Attention block.

1.1.3. Latent Variable Block

The primary distinction between the theoretical bound and the practical codec lies in their prior and posterior distributions. Specifically, in the practical codec's latent variable block, the posterior branch predicts only μ_i . In contrast, the bound model's posterior branch predicts both μ_i and σ_i . Figure 3 provides a detailed view of the latent variable block. The state of the practical codec latent variable block during training is shown in Figure 3(a). Figures 3(c) and 3(d) depict the encoding and decoding processes, respectively, while Figure 3(b) displays the latent variable block for the bound model.

Posteriors. The posterior distribution of Z_i given x and $z < i$ in the practical model is:

$$q_i \triangleq U\left(\mu_i - \frac{1}{2}, \mu_i + \frac{1}{2}\right) \quad (9)$$

$$\Leftrightarrow q_i(z_i | z_{<i}, x) = \begin{cases} 1 & \text{for } |z_i - \mu_i| \leq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases},$$

where μ_i is the output of the posterior branch, depending on the image x and preceding latent variables $z_{<i}$. Once q_i is obtained, z_i is sampled as $z_i \leftarrow \mu_i + u$, where u is a random sample from $U(-\frac{1}{2}, \frac{1}{2})$ during training, and during testing, it is replaced with scalar quantization.

In the theoretical bound model, the posterior distribution is a conditional Gaussian:

$$q_i \triangleq \mathcal{N}(\mu_i, \sigma_i^2) \quad (10)$$

Priors. For the practical model, the prior distribution p_i is defined as a conditional Gaussian convolved with a uniform distribution:

$$p_i \triangleq \mathcal{N}(\hat{\mu}_i, \hat{\sigma}_i^2) * U\left(-\frac{1}{2}, \frac{1}{2}\right) \quad (11)$$

$$\Leftrightarrow p_i(z_i | z_{<i}) = \int_{z_i - \frac{1}{2}}^{z_i + \frac{1}{2}} \mathcal{N}(t; \hat{\mu}_i, \hat{\sigma}_i^2) dt,$$

where $\mathcal{N}(t; \hat{\mu}_i, \hat{\sigma}_i^2)$ represents the Gaussian probability density function evaluated at t , and t is an integration dummy variable. The mean $\hat{\mu}_i$ and standard deviation $\hat{\sigma}_i$ are predicted by the prior branch. The probability mass function (PMF) P_i is then defined as:

$$P_i(n) \triangleq p_i(\hat{\mu}_i + n | z_{<i}), n \in \mathbb{Z}. \quad (12)$$

which is used for the entropy coding/decoding of z_i .

The theoretical bound model also employs a conditional Gaussian as the prior distribution:

$$p_i \triangleq \mathcal{N}(\hat{\mu}_i, \hat{\sigma}_i^2) \quad (13)$$

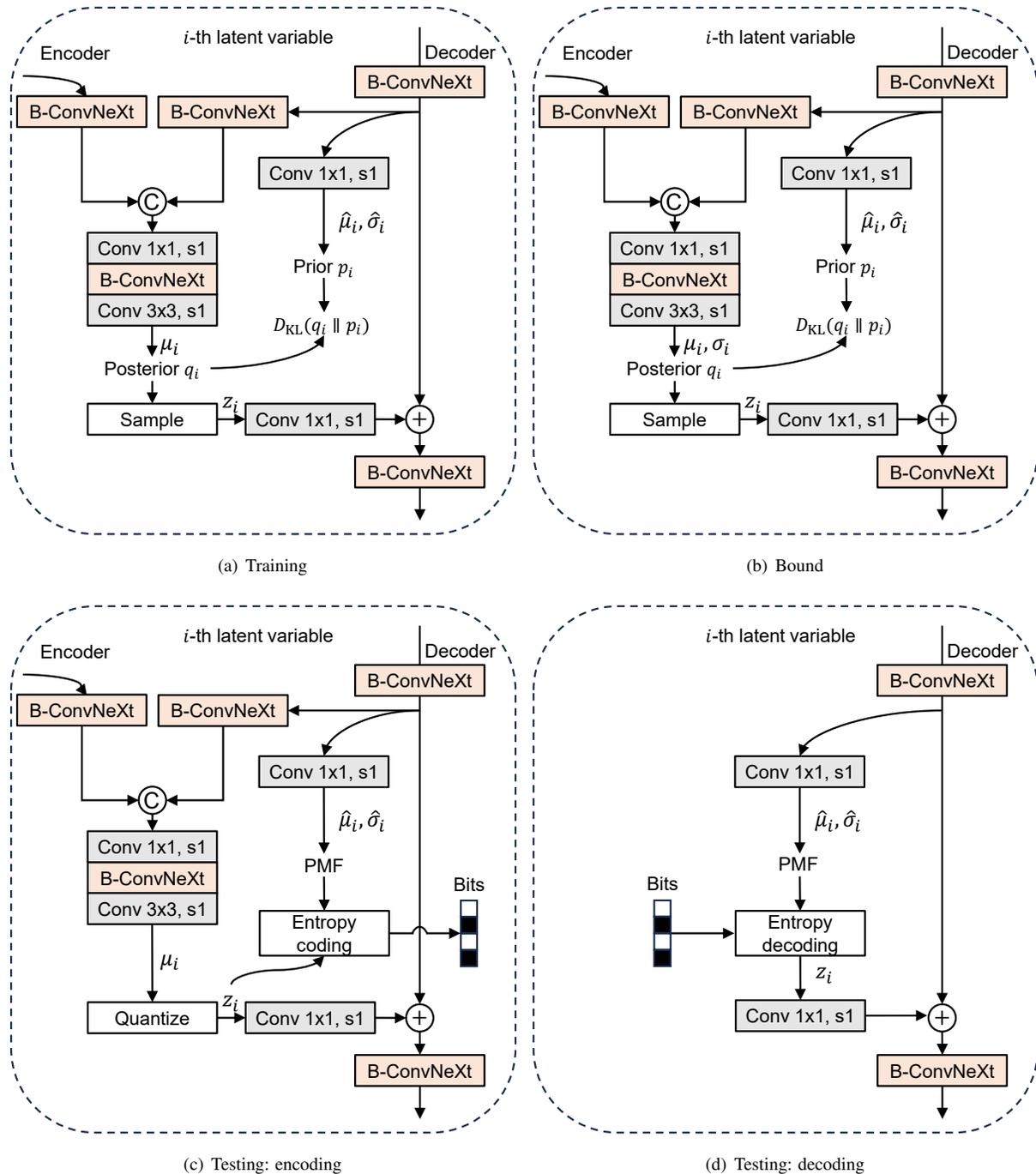


Fig. 3. Detailed Illustration of the Latent Variable Block.

1.2. Codecs Implementations

1.2.1. Learning-based Codecs

We list the implementations of the learning-based image codecs that we used for comparison in Table 1.

Table 1. Learning-based Codecs Implementations.

Method	Implementation
M&S Hyperprior [17]	github.com/InterDigitalInc/CompressAI
Cheng2020 [3]	github.com/InterDigitalInc/CompressAI
STF [1]	github.com/Googolxx/STF
TCM-S [5]	github.com/jmliu206/LIC_TCM
QARV [4]	github.com/duanzhihao/lossy-vaе

1.2.2. VVC Codec

We use VTM-18.0¹, the reference software for VVC, as the anchor. When testing VTM-18.0, we use OpenCV² to convert the image to YUV444 format and then use the following command line to test it. The output image is in YUV 4:4:4 format and converted back to RGB space. The final PSNR is computed between the final RGB image and the original RGB image.

```

1 EncoderApp.exe
2     -c encoder_intra_vtm.cfg
3     -o output.yuv
4     -q qp
5     -wdt image width
6     -hgt image height
7     -i input.yuv
8     -f 1
9     -fr 1
10    -fs 0
11    -b output.bin
12    --InputChromaFormat=444

```

1.3. Detailed Training Settings

We list detailed training information in Table 2, including data augmentation, hyperparameters, and training devices. We use different settings for the main experiment and the ablation experiments. In the main experiment, we train our model until convergence, which requires around 10 days of training on a dual-GPU machine. For ablation study experiments, we train our model with a shorter training period (500k iterations instead of 2M iterations) to reduce training costs.

1.4. RD curves Magnified

RD curves are magnified in Figures 4, 5, and 6 for the convenience of observation.

¹https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM/-/tree/VTM-18.0?ref_type=tags

²<https://github.com/opencv/opencv/tree/4.8.0>

Table 2. Training Hyperparameters.

	Main model	Ablation study
Training set	COCO 2017 train	COCO 2017 train
# images	118,287	118,287
Image size	Around 640x420	Around 640x420
Data augment.	Crop, h-flip	Crop, h-flip
Train input size	256x256	256x256
Optimizer	Adam	Adam
Learning rate	2×10^{-4}	2×10^{-4}
LR schedule	Constant + cosine	Constant
Batch size	32	32
# iterations	2M	500K
# images seen	64M	16M
Gradient clip	2.0	2.0
EMA	0.9999	0.9999
GPUs	2 × RTX 3090	1 × Quadro 6000
Time	260h	87h

2. REFERENCES

- [27] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen, “Conditional Positional Encodings for Vision Transformers,” *arXiv preprint arXiv:2102.10882*, 2021.

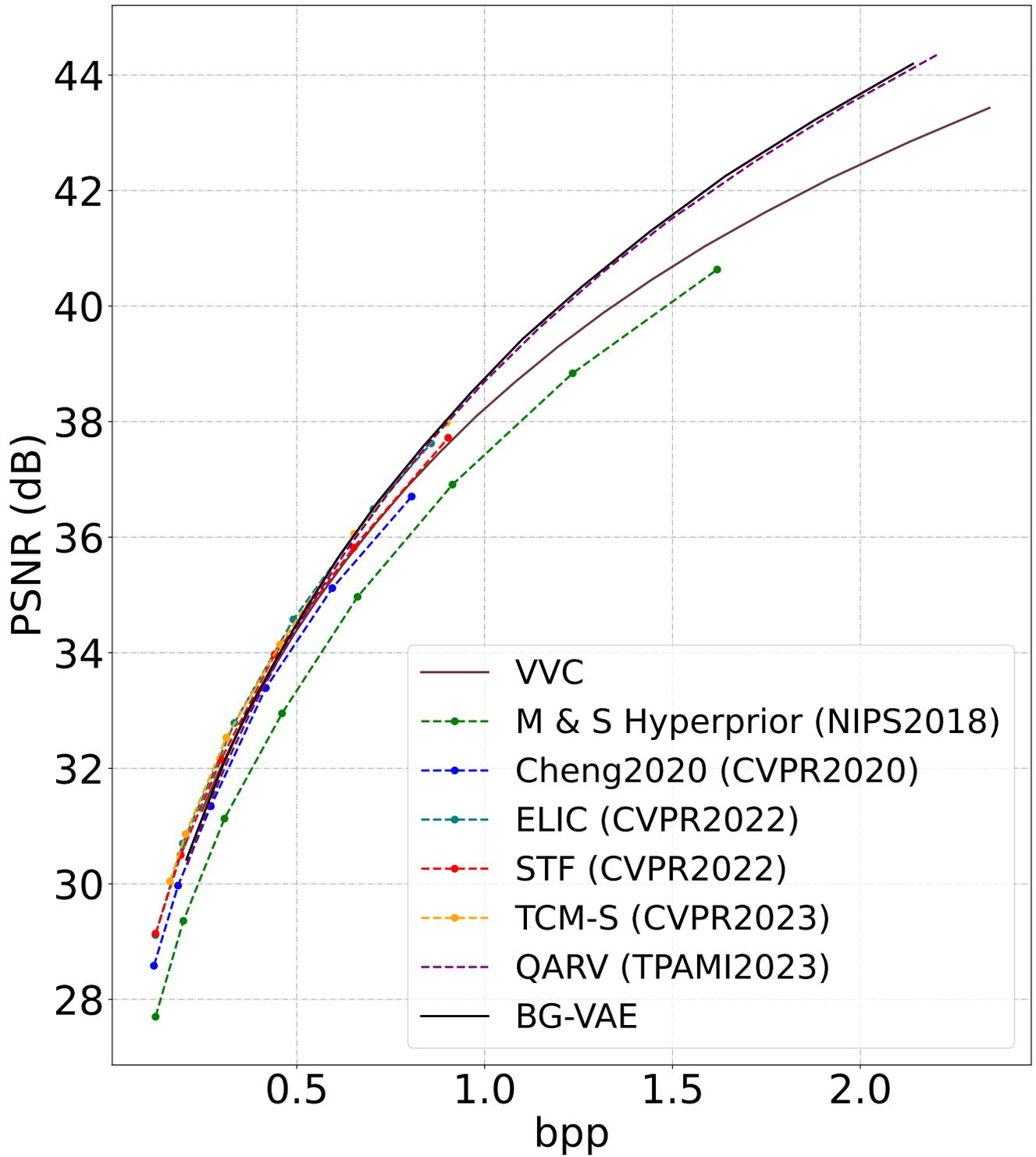


Fig. 4. RD Curves on Kodak Dataset.

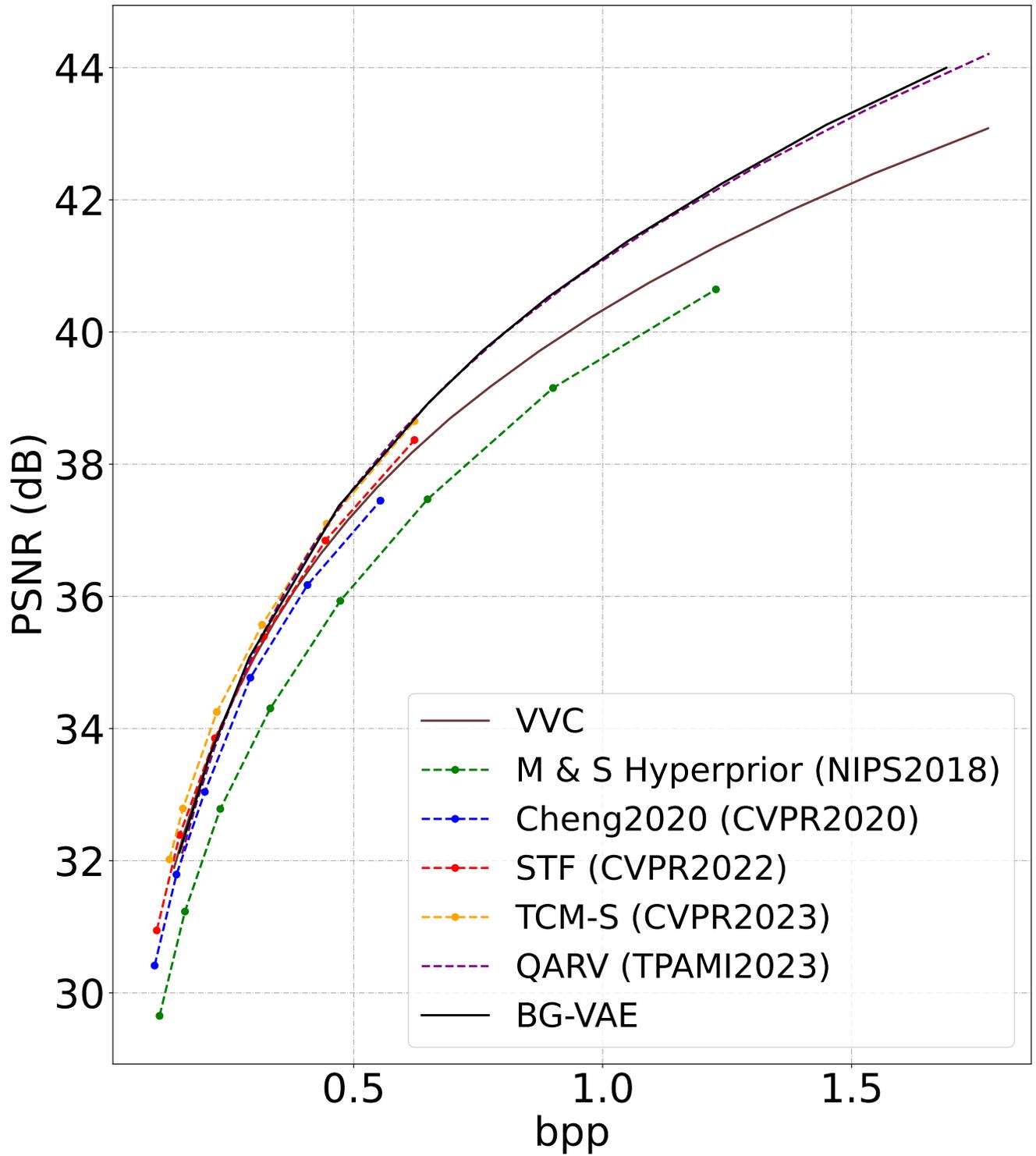


Fig. 5. RD Curves on Tecnick Dataset.

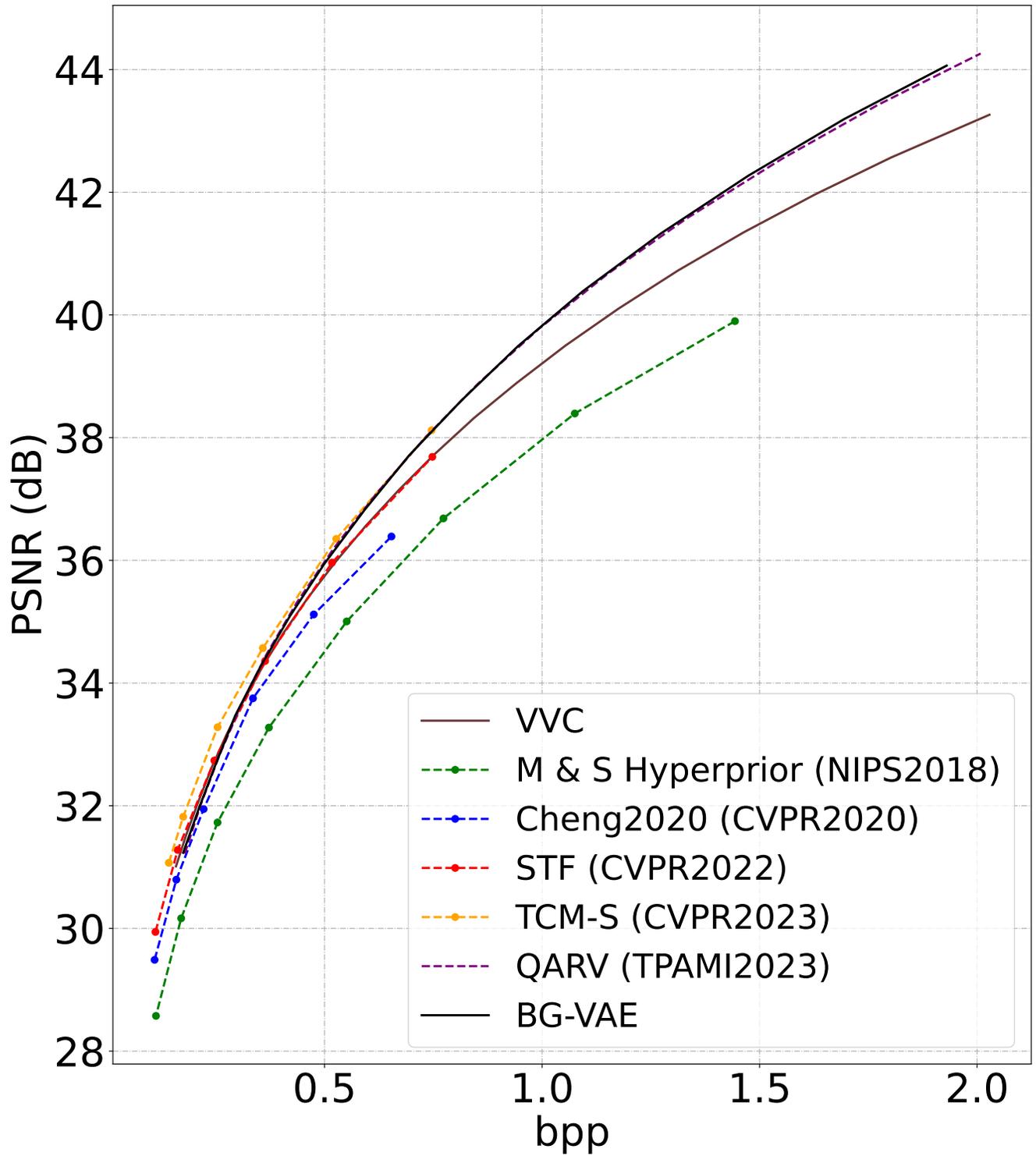


Fig. 6. RD Curves on CLIC2022 Dataset.