

ConvNeXt-ChARM: ConvNeXt-based Transform for Efficient Neural Image Compression

Ahmed Ghorbel

Wassim Hamidouche

Luce Morin

Abstract—Over the last few years, neural image compression has gained wide attention from research and industry, yielding promising end-to-end deep neural codecs outperforming their conventional counterparts in rate-distortion performance. Despite significant advancement, current methods, including attention-based transform coding, still need to be improved in reducing the coding rate while preserving the reconstruction fidelity, especially in non-homogeneous textured image areas. Those models also require more parameters and a higher decoding time. To tackle the above challenges, we propose ConvNeXt-ChARM, an efficient ConvNeXt-based transform coding framework, paired with a compute-efficient channel-wise auto-regressive prior to capturing both global and local contexts from the hyper and quantized latent representations. The proposed architecture can be optimized end-to-end to fully exploit the context information and extract compact latent representation while reconstructing higher-quality images. Experimental results on four widely-used datasets showed that ConvNeXt-ChARM brings consistent and significant BD-rate (PSNR) reductions estimated on average to 5.24% and 1.22% over the versatile video coding (VVC) reference encoder (VTM-18.0) and the state-of-the-art learned image compression method SwinT-ChARM, respectively. Moreover, we provide model scaling studies to verify the computational efficiency of our approach and conduct several objective and subjective analyses to bring to the fore the performance gap between the next generation ConvNet, namely ConvNeXt, and Swin Transformer. All materials, including the source code of SwinT-ChARM, will be made publicly accessible upon acceptance for reproducible research.

I. INTRODUCTION

Visual information is crucial in human development, communication, and engagement, and its compression is necessary for effective storage and transmission over constrained wireless/wireline channels. Thus, thinking about new lossy image compression approaches is a goldmine for scientific research. The goal is to reduce an image file size by permanently removing less critical information, particularly redundant data and high frequencies, to obtain the most compact bit-stream representation while preserving a certain level of visual fidelity. Nevertheless, the high compress rate and low distortion are fundamentally opposing objectives involving optimizing the rate-distortion tradeoff.

Conventional image and video compression standards including JPEG [1], JPEG2000 [2], H.265/high-efficiency video coding (HEVC) [3], and H.266/VVC [4], rely on hand-crafted

Ahmed Ghorbel and Luce Morin were with Univ. Rennes, INSA Rennes, CNRS, IETR - UMR 6164, Rennes, France, e-mail: (Ahmed.Ghorbel; Luce.Morin)@insa-rennes.fr.

Wassim Hamidouche was with Technology Innovation Institute, Masdar City, P.O Box 9639, Abu Dhabi, UAE, e-mail: Wassim.Hamidouche@tii.ae.

This work has been supported by Région Bretagne and Rennes Ville et Métropole under the DEEPTec project.

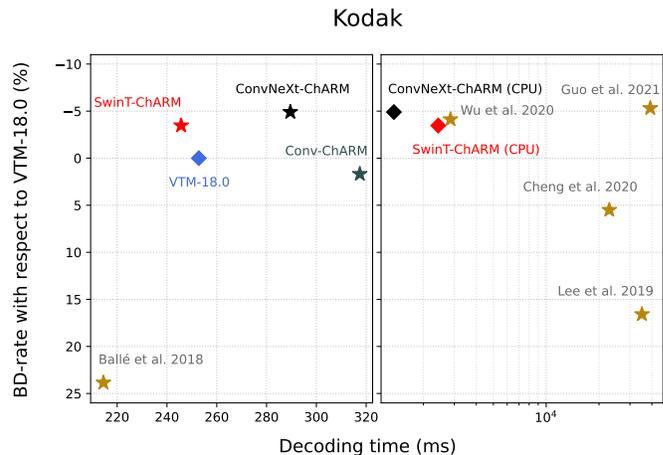


Fig. 1. BD-rate (%) versus decoding time (ms) on the Kodak dataset. Left-top is better. Star and diamond markers refer to decoding on GPU and CPU, respectively.

creativity to present module-based encoder/decoder block diagram. In addition, these codecs employ intra-prediction, fixed transform matrices, quantization, context-adaptive arithmetic coders, and various in-loop filters to reduce spatial and statistical redundancies, and alleviate coding artifacts. However, it has taken several years to standardize a conventional codec. Moreover, existing image compression standards are not anticipated to be an ideal and global solution for all types of image content due to the rapid development of new image formats and the growth of high-resolution mobile devices.

Lossy image compression consists of three modular parts: transform, quantization, and entropy coding. Each of these components can be represented as follows: i) autoencoders as flexible nonlinear transforms where the encoder (i.e., analysis transform) extracts latent representation from an input image and the decoder (i.e., synthesis transform) reconstructs the image from the decoded latent, ii) various differentiable quantization approaches which encode the latent into bitstream through arithmetic coding algorithms, iii) deep generative models as potent learnable entropy models estimating the conditional probability distribution of the latent to reduce the rate. Moreover, these three components can be optimized with end-to-end training by reducing the joint loss of the distortion between the original image and its reconstruction and the rate needed to transmit the bitstream of latent representation.

Thanks to recent advances in deep learning, we have seen many works exploring the potential of artificial neural networks (ANNs) to form various learned image and video

compression frameworks. Over the past two years, the performance of neural compression has steadily improved thanks to the prior line of study, reaching or outperforming state-of-the-art conventional codecs. Some previous works use local context [5]–[7], or additional side information [8]–[10] to capture short-range spatial dependencies, and others use non-local mechanism [11]–[14] as long-range spatial dependencies. Recently, Toderici *et al.* [15] proposed a generative compression method achieving high-quality reconstructions, Minnen *et al.* [16] introduced channel-conditioning and latent residual prediction taking advantage of an entropy-constrained model that uses both forward and backward adaptations, and Zhu *et al.* [17] replaced all convolutions in the channel-wise autoregressive model (ChARM) prior approach [16] with Swin Transformer [18] blocks, Zou *et al.* [19] combined the local-aware attention mechanism with the global-related feature learning and proposed a window-based attention module, Koyuncu *et al.* [20] proposed a Transformer-based context model, which generalizes the standard attention mechanism to spatio-channel attention, Zhu *et al.* [21] proposed a probabilistic vector quantization with cascaded estimation under a multi-codebooks structure, Kim *et al.* [22] exploited the joint global and local hyperpriors information in a content-dependent manner using an attention mechanism, and He *et al.* [23] adopted stacked residual blocks as nonlinear transform and multi-dimension entropy estimation model.

One of the main challenges of learned transform coding is the ability to identify the crucial information necessary for the reconstruction, knowing that information overlooked during encoding is usually lost and unrecoverable for decoding. Another main challenge is the tradeoff between performance and decoding speed. While the existing approaches improve the transform and entropy coding accuracy, they remain limited by the higher decoding runtime and excessive model complexity leading to an ineffective real-world use. Finally, we found that attention-based networks taking advantage of attention mechanisms to capture global dependencies, such as Swin Transformer [18], have over-smoothed and contain undesirable artifacts at low bitrates. Furthermore, the global semantic information in image compression is less effective than in other computer vision tasks [19].

In this paper, we propose a nonlinear transform built on ConvNeXt blocks with additional down and up sampling layers and paired with a ChARM prior, namely ConvNeXt-ChARM. Recently proposed in [24], ConvNeXt is defined as a modernized ResNet architecture toward the design of a vision Transformer, which competes favorably with Transformers in terms of efficiency, achieving state-of-the-art on ImageNet classification task [25] and outperforming Swin Transformer on COCO detection [26] and ADE20K segmentation [27] challenges while maintaining the maturity and simplicity of convolutional neural networks (ConvNets) [24]. The contributions of this paper are summarized as follows:

- We propose a learned image compression model that leverages a stack of ConvNeXt blocks with down and up-sampling layers for extracting contextualized and non-

linear information for effective latent decorrelation. We maintain the convolution strengths like sliding window strategy for computations sharing, translation equivariance as a built-in inductive bias, and the local nature of features, which are intrinsic to providing a better spatial representation.

- We apply ConvNeXt-based transform coding layers for generating and decoding both latent and hyper-latent to consciously and subtly balance the importance of feature compression through the end-to-end learning framework.
- We conduct experiments on four widely-used evaluation datasets to explore possible coding gain sources and demonstrate the effectiveness of ConvNeXt-ChARM. In addition, we carried out a model scaling analysis to compare the complexity of ConvNeXt and Swin Transformer.

Extensive experiments validate that the proposed ConvNeXt-ChARM achieves state-of-the-art compression performance, as illustrated in Figure 1, outperforming conventional and learned image compression methods in the tradeoff between coding efficiency and decoder complexity.

The rest of this paper is organized as follows. Section II presents our overall framework along with a detailed description of the proposed architecture. Next, we dedicate Section III to describe and analyze the experimental results. Finally, Section IV concludes the paper.

II. PROPOSED CONVNEXT-CHARM MODEL

A. Problem Formulation

The objective of learned image compression is to minimize the distortion between the original image and its reconstruction under a specific distortion-controlling hyper-parameter. Assuming an input image \mathbf{x} , the analysis transform g_a , with parameter ϕ_g , removes the image spatial redundancies and generates the latent representation \mathbf{y} . Then, this latent is quantized to the discrete code $\hat{\mathbf{y}}$ using the quantization operator $\lceil \cdot \rceil$, from which a synthesis transform g_s , with parameter θ_g , reconstructs the image denoted by $\hat{\mathbf{x}}$. The overall process can be formulated as follows:

$$\begin{aligned} \mathbf{y} &= g_a(\mathbf{x} \mid \phi_g), \\ \hat{\mathbf{y}} &= \lceil \mathbf{y} \rceil, \\ \hat{\mathbf{x}} &= g_s(\hat{\mathbf{y}} \mid \theta_g). \end{aligned} \quad (1)$$

A hyperprior model composed of a hyper-analysis and hyper-synthesis transforms (h_a, h_s) with parameters (ϕ_h, θ_h) is usually used to reduce the statistical redundancy among latent variables. In particular, this hyperprior model assigns a few extra bits as side information to transmit some spatial structure information and helps to learn an accurate entropy model. The hyperprior generation can be summarized as follows:

$$\begin{aligned} \mathbf{z} &= h_a(\mathbf{y} \mid \phi_h), \\ \hat{\mathbf{z}} &= \lceil \mathbf{z} \rceil, \\ p_{\hat{\mathbf{y}}|\hat{\mathbf{z}}}(\hat{\mathbf{y}} \mid \hat{\mathbf{z}}) &\leftarrow h_s(\hat{\mathbf{z}} \mid \theta_h). \end{aligned} \quad (2)$$

Transform and quantization introduce a distortion $D = MSE(\mathbf{x}, \hat{\mathbf{x}})$, for mean squared error (MSE) optimization that

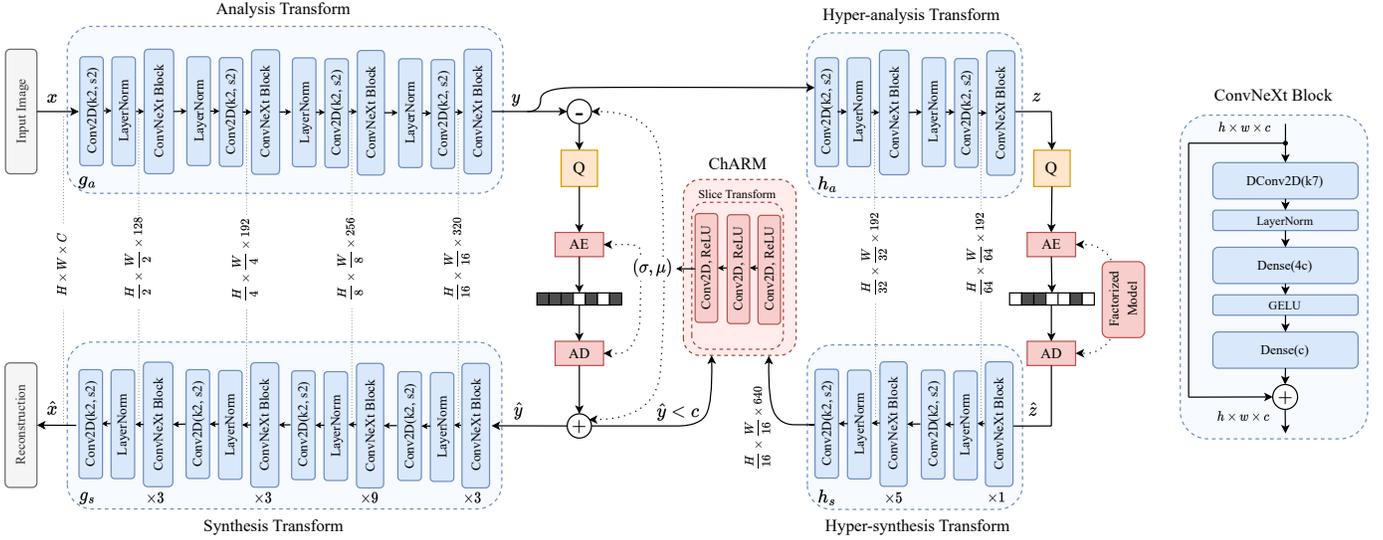


Fig. 2. Overall ConvNeXt-ChARM Framework. We illustrate the image compression diagram of our ConvNeXt-ChARM with hyperprior and channel-wise auto-regressive context model. We also present the ConvNeXt block used in both transform and hyper-transform coding for an end-to-end feature aggregation.

measures the reconstruction quality with an estimated bitrate R , corresponding to the expected rate of the quantized latents and hyper-latents, as described below:

$$R = \mathbb{E} \left[-\log_2(p_{\hat{y}|\hat{z}}(\hat{y}|\hat{z})) - \log_2(p_{\hat{z}}(\hat{z})) \right]. \quad (3)$$

Representing (g_a, g_s) , (h_a, h_s) , and entropy model by deep neural networks (DNNs) enables jointly optimizing the end-to-end model by minimizing the rate-distortion tradeoff \mathcal{L} , giving a rate-controlling hyper-parameter λ . This optimization problem can be presented as follows:

$$\begin{aligned} \mathcal{L} &= R + \lambda D, \\ &= \underbrace{\mathbb{H}(\hat{y}) + \mathbb{H}(\hat{z})}_R + \lambda MSE(x, \hat{x}), \end{aligned} \quad (4)$$

where \mathbb{H} stands for the entropy.

B. ConvNeXt-ChARM network architecture

To better parameterize the distributions of the quantized latent features with a more accurate and flexible entropy model, we adopted the ChARM prior approach proposed in [16] to build an efficient ConvNeXt-based learning image compression model with strong compression performance. As shown in Figure 2, the analysis/synthesis transform (g_a, g_s) of our design consists of a combination of down and up-sampling blocks and ConvNeXt encoding/decoding blocks [24], respectively. Down and up-sampling blocks are performed using Conv2D and Normalisation layers sequentially. The architectures for hyper-transforms (h_a, h_s) are similar to (g_a, g_s) with different stages and configurations.

C. ConvNeXt design description

Globally, ConvNeXt incorporates a series of architectural choices from a Swin Transformer while maintaining the network's simplicity as a standard ConvNet without introducing any attention-based modules. These design

decisions can be summarized as follows: macro design, ResNeXt's grouped convolution, inverted bottleneck, large kernel size, and various layer-wise micro designs. In Figure 2, we illustrates the ConvNeXt block, where the DConv2D(.) refers for the a depthwise 2D convolution, LayerNorm for the layer normalization, Dense(.) for the densely-connected NN layer, and GELU for the activation function.

Macro design: The stage compute ratio is adjusted from (3, 4, 6, 3) in ResNet-50 to (3, 3, 9, 3), which also aligns the FLOPs with Swin-T. In addition, the ResNet-style stem cell is replaced with a patchify layer implemented using a 2×2 , stride two non-overlapping convolutional layers with an additional normalization layer to help stabilize the training. In ConvNeXt-ChARM diagram, we adopted the (3, 3, 9, 3) and (5, 1) as stage compute ratios for transforms and hyper-transforms, respectively.

Depthwise convolution: The ConvNeXt block uses a depthwise convolution, a special case of grouped convolution used in ResNeXt [28], where the number of groups is equal to the considered channels. This is similar to the weighted sum operation in self-attention, which operates by mixing information only in the spatial dimension.

Inverted bottleneck: Similar to Transformers, ConvNeXt is designed with an inverted bottleneck block, where the hidden dimension of the residual block is four times wider than the input dimension. As illustrated in the ConvNeXt block Figure 2, the first dense layer is 4 times wider then the second one.

large kernel: One of the most distinguishing aspects of Swin Transformers is their local window in the self-attention

block. The information is propagated across windows, which enables each layer to have a global receptive field. The local window is at least 7×7 sized, which is still more extensive than the 3×3 ResNeXt kernel size. Therefore, ConvNeXt adopted large kernel-sized convolutions by using a 7×7 depthwise 2D convolution layer in each block. This allows our ConvNeXt-ChARM model to capture global contexts in both latents and hyper-latents, which are intrinsic to providing a better spatial representation.

Micro design: In ConvNeXt’s micro-design, several per-layer enhancements are applied in each block, by using: a single Gaussian error linear unit (GELU) activation function (instead of numerous ReLU), using a single LayerNorm as normalization choice (instead of numerous BatchNorm), and using separate down-sampling layers between stages.

III. RESULTS

First, we briefly describe used datasets with the implementation details. Then, we assess the compression efficiency of our method with a rate-distortion comparison and compute the average bitrate savings on four commonly-used evaluation datasets. We further elaborate a model scaling and complexity study to consistently examine the effectiveness of our proposed method against pioneering ones.

A. Experimental Setup

Datasets. The training set of the CLIC2020 dataset is used to train the proposed ConvNeXt-ChARM model. This dataset contains a mix of professional and user-generated content images in RGB color and grayscale formats. We evaluate image compression models on four datasets, including Kodak [29], Tecnick [29], JPEG-AI [29], and the testing set of CLIC21 [29]. For a fair comparison, all images are cropped to the highest possible multiples of 256 to avoid padding for neural codecs.

Implementation details. We implemented all models in TensorFlow using tensorflow compression (TFC) library [30], and the experimental study was carried out on an RTX 5000 Ti GPU. All models were trained on the same CLIC2020 training set with 3.5M steps using the ADAM optimizer with parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The initial learning rate is set to 10^{-4} and drops to 10^{-5} for another 100k iterations, and $L = R + \lambda D$ as loss function. The MSE is used as the distortion metric in RGB color space. Each batch contains eight random 256×256 crops from training images. To cover a wide range of rate and distortion, for our proposed method, we trained five models with $\lambda \in \{0.006, 0.009, 0.020, 0.050, 0.150\}$. Regarding the evaluation on CPU, we used an Intel(R) Xeon(R) W-2145 @ 3.70GHz.

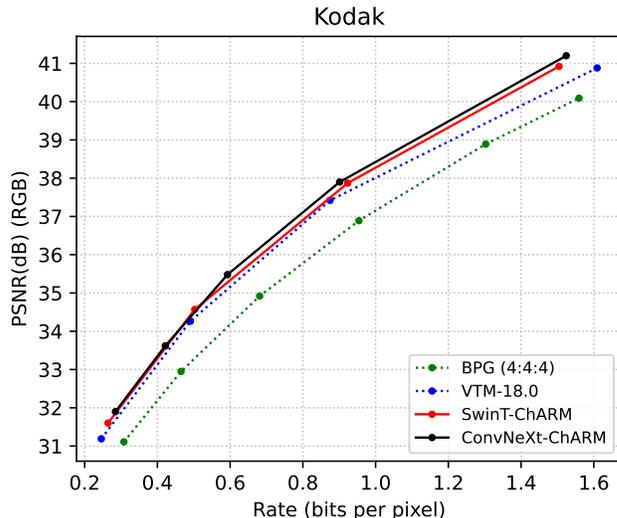


Fig. 3. Rate-distortion comparison on Kodak dataset.

Baselines.¹ We compare our approach with the state-of-art neural compression method SwinT-ChARM proposed by Zhu *et al.* [17], and non-neural compression methods, including better portable graphics (BPG)(4:4:4), and the most up-to-date VVC official Test Model VTM-18.0 in All-Intra profile configuration.

B. Rate-Distortion coding performance

To demonstrate the compression efficiency of our proposed approach, we visualize the rate-distortion curves of our model and the baselines on each of the considered datasets. Considering the Kodak dataset, Figure 3 shows that our ConvNeXt-ChARM outperforms the state-of-the-art learned approach SwinT-ChARM, as well as the BPG(4:4:4) and VTM-18.0 traditional codecs in terms of PSNR. Regarding rate savings over VTM-18.0, SwinT-ChARM has more compression abilities only for low PSNR values. Our model can be generalized to high resolution image datasets (Tecnick, JPEG-AI, and CLIC21), and can still outperform existing traditional and the learned image compression method SwinT-ChARM in terms of PSNR. Besides the rate-distortion curves, we also evaluate different models using Bjontegaard’s metric [31], which computes the average bitrate savings (%) between two rate-distortion curves. In Table I, we summarize the BD-rate of image codecs across all four datasets compared to the VTM-18.0 as the anchor. On average, ConvNeXt-ChARM is able to achieve 5.24% rate reduction compared to VTM-18.0 and 1.22% relative gain from SwinT-ChARM. Figure 1 shows the BD-rate (with VTM-18.0 as an anchor) versus the decoding time of various approaches on the Kodak dataset. It can be seen from the figure that our ConvNeXt-ChARM achieves a good tradeoff between BD-rate performance and decoding time.

¹For a fair comparison, we only considered SwinT-ChARM [17] from the state-of-the-art models [17], [19]–[23], due to the technical feasibility of models training and evaluation under the same conditions and in an adequate time.

TABLE I
BD-RATE \downarrow PERFORMANCE OF BPG (4:4:4), SWIN-T-CHARM, AND CONVNEXT-CHARM COMPARED TO THE VTM-18.0 FOR THE FOUR CONSIDERED DATASETS.

Dataset	BPG444	SwinT-ChARM	ConvNeXt-ChARM
Kodak	20.73%	-3.47%	-4.90%
Tecnick	27.03%	-6.52%	-7.56%
JPEG-AI	28.14%	-0.23%	-1.17%
CLIC21	26.54%	-5.86%	-7.36%
Average	25.61%	-4.02%	-5.24%

TABLE II
IMAGE CODEC COMPLEXITY. WE CALCULATED THE AVERAGE DECODING TIME ACROSS 7000 IMAGES AT 256×256 RESOLUTION, ENCODED AT 0.6 BPP. THE BEST SCORE IS HIGHLIGHTED IN BOLD.

Image Codec	Latency(ms) \downarrow		GFLOPs \downarrow	#params(M) \downarrow
	GPU	CPU		
Conv-ChARM	124.32	967.43	117	123.84
SwinT-ChARM	102.45	1088.16	122	127.78
Ours	122.70	834.42	119	122.33

C. Models Scaling Study

We evaluated the decoding complexity of the three considered image codecs by averaging decoding time across 7000 images at 256×256 resolution, encoded at 0.6 bpp. We present the image codec complexity in Table II, including decoding time on GPU and CPU, floating point operations per second (GFLOPs), the memory required by model weights, and the total model parameters. The models run with Tensorflow 2.8 on a workstation with one RTX 5000 Ti GPU. The Conv-ChARM model refers to the Minnen *et al.* [16] architecture with a latent depth of 320 and a hyperprior depth of 192, and can be considered as ablation of our model without ConvNeXt blocks. We maintained the same slice transform configuration of the ChARM for the three considered models. The total decoding time of SwinT-ChARM decoder is less than ConvNets-based decoder on GPU but is the highest on CPU. Our ConvNeXt-ChARM is lighter than the Conv-ChARM in terms of the number of parameters, which proves the ConvNeXt block’s well-engineered design. Compared with SwinT-ChARM, our ConvNeXt-ChARM shows lower complexity, requiring lower training time with less memory consumption. In addition, Figure 4 shows that our method is in an interesting area, achieving a good tradeoff between BD-rate score on Kodak, total model parameters, and MFLOPs per pixel, highlighting an efficient and hardware-friendly compression model.

D. Comparison with SwinT-ChARM

ConvNeXt-ChARM achieves good rate-distortion performance while significantly reducing the latency, which is potentially helpful to conduct, with further optimizations, high-quality real-time visual data transmission, as recently proposed in the first software-based neural video decoder running HD resolution video in real-time on a commercial smartphone [32]. Since fewer works attempt to explicitly compare Swin

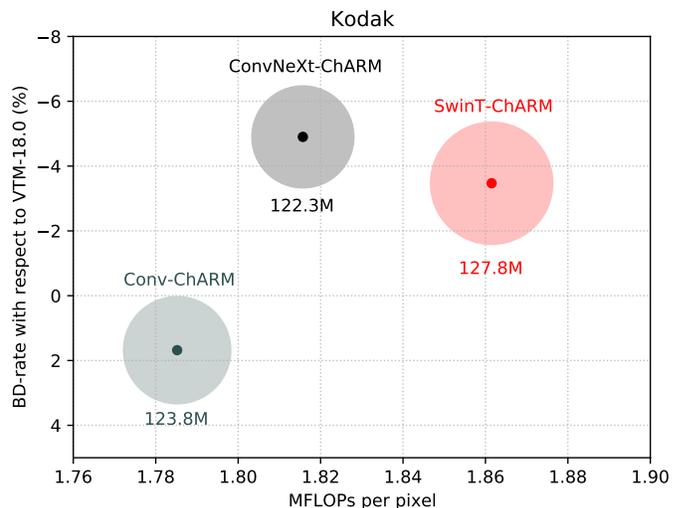


Fig. 4. Model size scaling. BD-Rate versus MFLOPs per pixel for our model ConvNeXt-ChARM compared to Conv-ChARM and SwinT-ChARM (for both encoding and decoding).

Transformer and ConvNet-based blocks, here, we compare our ConvNeXt-ChARM with SwinT-ChARM under the same conditions and configurations. We found that a well-designed ConvNet, without any additional attention modules, can outperform the highly coveted Swin Transformer in learned transform coding in terms of BD-rate, with more visually pleasing reconstructions and comparable decoding latency. In addition, ConvNeXt-ChARM maintains the efficiency and maturity of standard ConvNets and the fully-convolutional nature for both training and inference. There is no doubt that Transformers are excellent architectures with enormous potential for the future of various computer vision applications. However, their vast hunger for data and computational resources [33] poses a big challenge for the computer vision community. Taking SwinT-ChARM as an example, it needs, on average, $\times 1.33$ more time than ConvNeXt-ChARM, to train on the same number of epochs.

IV. CONCLUSION

In this work, we reconcile compression efficiency with ConvNeXt-based transform coding paired with a ChARM prior and propose an up-and-coming learned image compression model ConvNeXt-ChARM. Furthermore, we inherit the advantages of pure ConvNets in the proposed method to improve both efficiency and effectiveness. The experimental results, conducted on four datasets, showed that our approach outperforms previously learned and conventional image compression methods, creating a new state-of-the-art rate-distortion performance with a significant decoding runtime decrease. Future work will further investigate efficient low-complexity entropy coding approaches to further enhance decoding latency. With the development of GPU chip technology and the further optimization of engineering, learning-based codecs will be the future of coding, achieving better

compression efficiency when compared with traditional codecs and aiming to bridge the gap to a real-time operation. We hope our study will challenge certain accepted notions and prompt people to reconsider the significance of convolutions in computer vision.

REFERENCES

- [1] Ricardo Monteiro, Luis Lucas, Caroline Conti, Paulo Nunes, Nuno Rodrigues, Sérgio Faria, Carla Pagliari, Eduardo Da Silva, and Luís Soares, "Light field hevc-based image coding using locally linear embedding and self-similarity compensated prediction," in *2016 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2016, pp. 1–4.
- [2] Majid Rabbani and Rajan Joshi, "An overview of the jpeg 2000 still image compression standard," *Signal processing: Image communication*, vol. 17, no. 1, pp. 3–48, 2002.
- [3] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand, "Overview of the high efficiency video coding (hevc) standard," *IEEE Transactions on circuits and systems for video technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [4] Gary Sullivan, "Versatile video coding (vvc) arrives," in *2020 IEEE International Conference on Visual Communications and Image Processing (VCIP)*. IEEE, 2020, pp. 1–1.
- [5] David Minnen, Johannes Ballé, and George D Toderici, "Joint autoregressive and hierarchical priors for learned image compression," *Advances in neural information processing systems*, vol. 31, 2018.
- [6] Jooyoung Lee, Seunghyun Cho, Seyoon Jeong, Hyoungjin Kwon, Hyunsuk Ko, Hui Yong Kim, and Jin Soo Choi, "Extended end-to-end optimized image compression method based on a context-adaptive entropy model," in *CVPR Workshops*, 2019, p. 0.
- [7] Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool, "Conditional probability models for deep image compression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4394–4402.
- [8] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston, "Variational image compression with a scale hyperprior," *arXiv preprint arXiv:1802.01436*, 2018.
- [9] Yueyu Hu, Wenhan Yang, and Jiaying Liu, "Coarse-to-fine hyper-prior modeling for learned image compression," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, vol. 34, pp. 11013–11020.
- [10] David Minnen, George Toderici, Saurabh Singh, Sung Jin Hwang, and Michele Covell, "Image-dependent local entropy models for learned image compression," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 430–434.
- [11] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto, "Learned image compression with discretized gaussian mixture likelihoods and attention modules," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7939–7948.
- [12] Mu Li, Kai Zhang, Jinxing Li, Wangmeng Zuo, Radu Timofte, and David Zhang, "Learning context-based nonlocal entropy modeling for image compression," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [13] Yichen Qian, Zhiyu Tan, Xiuyu Sun, Ming Lin, Dongyang Li, Zhenhong Sun, Hao Li, and Rong Jin, "Learning accurate entropy model with global reference for image compression," *arXiv preprint arXiv:2010.08321*, 2020.
- [14] Tong Chen, Haojie Liu, Zhan Ma, Qiu Shen, Xun Cao, and Yao Wang, "End-to-end learnt image compression via non-local attention optimization and improved context modeling," *IEEE Transactions on Image Processing*, vol. 30, pp. 3179–3191, 2021.
- [15] George Dan Toderici, Fabian Julius Mentzer, Eirikur Thor Agustsson, and Michael Tobias Tschannen, "High-fidelity generative image compression," June 2 2022, US Patent App. 17/107,684.
- [16] David Minnen and Saurabh Singh, "Channel-wise autoregressive entropy models for learned image compression," in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 3339–3343.
- [17] Yin hao Zhu, Yang Yang, and Taco Cohen, "Transformer-based transform coding," in *International Conference on Learning Representations*, 2021.
- [18] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [19] Renjie Zou, Chunfeng Song, and Zhaoxiang Zhang, "The devil is in the details: Window-based attention for image compression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17492–17501.
- [20] A Burakhan Koyuncu, Han Gao, Atanas Boev, Georgii Gaikov, Elena Alshina, and Eckehard Steinbach, "Contextformer: A transformer with spatio-channel attention for context modeling in learned image compression," in *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIX*. Springer, 2022, pp. 447–463.
- [21] Xiaosu Zhu, Jingkuan Song, Lianli Gao, Feng Zheng, and Heng Tao Shen, "Unified multivariate gaussian mixture for efficient neural image compression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17612–17621.
- [22] Jun-Hyuk Kim, Byeongho Heo, and Jong-Seok Lee, "Joint global and local hierarchical priors for learned image compression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5992–6001.
- [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," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5718–5727.
- [24] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11976–11986.
- [25] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [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," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [27] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba, "Scene parsing through ade20k dataset," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 633–641.
- [28] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [29] Benchmark datasets, "kodak testing set: <http://r0k.us/graphics>, technic testing set: <https://testimages.org/>, jpeg-ai testing set: https://jpegai.github.io/test_images/, and clic21 testing set: <http://compression.cc/tasks/>," .
- [30] Johannes Ballé, Sung Jin Hwang, and Eirikur Agustsson, "TensorFlow Compression: Learned data compression," 2022.
- [31] Gisle Bjontegaard, "Calculation of average psnr differences between rd-curves," *VCEG-M33*, 2001.
- [32] Hoang Le, Liang Zhang, Amir Said, Guillaume Sautiere, Yang Yang, Pranav Shrestha, Fei Yin, Reza Pourreza, and Auke Wiggers, "Mobilecodec: neural inter-frame video compression on mobile devices," in *Proceedings of the 13th ACM Multimedia Systems Conference*, 2022, pp. 324–330.
- [33] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.