

AICT: AN ADAPTIVE IMAGE COMPRESSION TRANSFORMER

Ahmed Ghorbel¹ Wassim Hamidouche^{1,2} Luce Morin¹

¹ Univ Rennes, INSA Rennes, CNRS, IETR – UMR 6164, F-35000 Rennes, France

² Technology Innovation Institute P.O.Box: 9639, Masdar City Abu Dhabi, UAE

ABSTRACT

Motivated by the efficiency investigation of the Transformer-based transform coding framework, namely SwinT-ChARM, we propose to enhance the latter, as first, with a more straightforward yet effective Transformer-based channel-wise auto-regressive prior model, resulting in an absolute image compression transformer (ICT). Current methods that still rely on ConvNet-based entropy coding are limited in long-range modeling dependencies due to their local connectivity and an increasing number of architectural biases and priors. On the contrary, the proposed ICT can capture both global and local contexts from the latent representations and better parameterize the distribution of the quantized latents. Further, we leverage a learnable scaling module with a sandwich ConvNeXt-based pre/post-processor to accurately extract more compact latent representation while reconstructing higher-quality images. Extensive experimental results on benchmark datasets showed that the proposed adaptive image compression transformer (AICT) framework significantly improves the trade-off between coding efficiency and decoder complexity over the versatile video coding (VVC) reference encoder (VTM-18.0) and the neural codec SwinT-ChARM.

Index Terms— Neural Image Compression, Adaptive Resolution, Spatio-Channel Entropy Modeling, Self-attention, Transformer.

1. INTRODUCTION

Visual information is crucial in human development, communication, and engagement, and its compression is necessary for effective data storage and transmission. Thus, designing new lossy image compression algorithms is a goldmine for scientific research. The goal is to reduce an image file size by permanently removing less critical information, specifically redundant data and high frequencies, to obtain the most compact bit-stream representation while preserving a certain level of visual fidelity. Nevertheless, the high compression rate and low distortion are fundamentally opposed objectives involving optimizing the rate-distortion (RD) cost.

Conventional compression standards, including JPEG, JPEG2000, H.265/HEVC, and H.266/VVC, rely on hand-crafted creativity to present module-based encoder/decoder block diagram, i.e., Intra prediction, transform, quantization, arithmetic coding, and post-processing. Traditional coding algorithms have a lot of advantages, including mature technology with SW/HW-friendly implementations, low decoding complexity, and strong generalization on different contents. Nevertheless, all of them mainly rely on hand-crafted coding techniques; thus, it is quite challenging to directly optimize RD cost for all types of image content due to the rapid development of new image formats and the growth of high-resolution

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

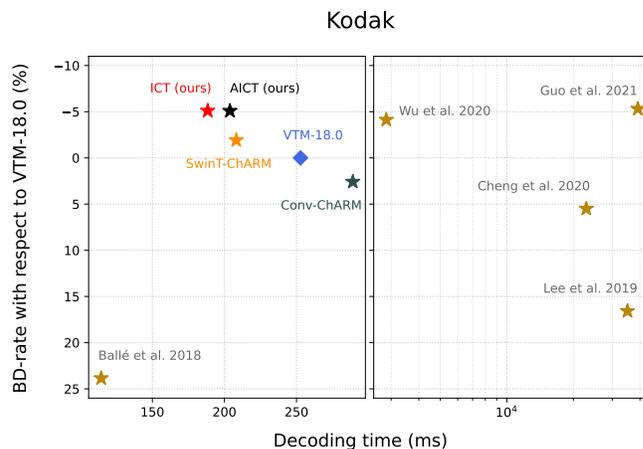


Fig. 1. BD-rate (%)↓ versus decoding time (ms)↓ on the Kodak dataset. Left-top is better. Star and diamond markers refer to decoding on GPU and CPU, respectively.

mobile devices. On the other hand, neural coding has gained wide attention from research and industry, yielding promising end-to-end neural image compression (NIC) solutions outperforming their conventional counterparts in coding efficiency. NIC leverages autoencoders (AEs) to carry out a non-linear coding from the signal domain to a compact representation. Such AE-based system consists of three modular parts: transform, quantization, and entropy coding, which can be trained in an end-to-end fashion to minimize the distortion between a source image and its reconstruction, and the rate needed to convey the latent representation bit-stream.

Since the early recurrent neural network (RNN)-based method [1] for lossy image compression, significant advancements have been made in integrating tailored modules for NIC. Previous works use local context [2–4], or additional side information [5–7] to capture short-range spatial dependencies, and others use non-local mechanisms [8–11] to model long-range spatial dependencies. Recently, Toderici *et al.* [12] proposed a generative compression method achieving high-quality reconstructions, Minnen *et al.* [13] introduced channel-conditioning taking advantage of an entropy-constrained model that uses both forward and backward adaptations, Zhu *et al.* [14] replaced the ConvNet-based transform coding in the Minnen *et al.* [13] approach with a Transformer-based one, Zou *et al.* [15] combined the local-aware attention mechanism with the global-related feature learning and proposed a window-based attention module, Koyuncu *et al.* [16] proposed a Transformer-based context model, which generalizes the standard attention mechanism to spatio-channel attention, Zhu *et al.* [17] proposed a probabilistic vector quantization with cascaded estimation under a multi-codebooks structure, Kim *et al.* [18] exploited the joint global and local hyper-

priors information in a content-dependent manner using an attention mechanism, and He *et al.* [19] adopted stacked residual blocks as nonlinear transform and multi-dimension entropy estimation model.

In order to improve image-level prediction while minimizing computation costs, learned sampling techniques have been developed for several vision tasks. Spatial transformer networks (STNs) [20] introduce a layer that estimates a parametrized affine, projective, and splines transformation from an input image to recover data distortions. Based on the latter, Chen *et al.* [21] proposed a straightforward learned downsampling module that can be jointly optimized with any neural compression kernels in an end-to-end fashion. Talebi *et al.* [22] jointly optimize pixel value interpolated at each fixed downsampling location for classification. Jin *et al.* [23] introduced a deformation module and a learnable downsampling operation, which can be optimized together with the given segmentation model.

One of the main challenges of NIC 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 trade-off between coding performance and decoding latency. While the existing approaches improve the transform and entropy coding accuracy, they still need to be improved by the higher decoding runtime and excessive model complexity leading to an ineffective real-world use. To cope with those challenges, we present in this paper three contributions summarized as follows:

- We propose the image compression transformer (ICT), a nonlinear transform coding and spatio-channel auto-regressive entropy coding. These modules are based on swin transformer (SwinT) blocks for effective latent decorrelation and a more flexible receptive field to adapt for contexts requiring short/long-range information.
- We further propose the adaptive image compression transformer (AICT) model that adopts a scale adaptation module as a sandwich processor to enhance compression efficiency. This module consists of a neural scaling network and ConvNeXt-based pre/post-processor to jointly optimize differentiable resizing layers and a content-dependent resize factor estimator.
- We conduct experiments on four widely-used benchmark datasets to explore possible coding gain sources and demonstrate the effectiveness of AICT. In addition, we carried out a model scaling analysis and an ablation study to substantiate our architectural decisions.

Extensive experiments reveal the impacts of the spatio-channel entropy coding, the sandwich scale adaptation component, and the joint global structure and local texture learned by the self-attention units through the nonlinear transform coding. These experiments validate that the proposed AICT model achieves compelling compression performance, as illustrated in Fig. 1, outperforming conventional and neural codecs in both coding efficiency and decoder complexity.

The rest of this paper is organized as follows. First, the proposed AICT framework is described in detail in Section 2. Next, we dedicate Section 3 to describe and analyze the experimental results. Finally, Section 4 concludes the paper.

2. PROPOSED AICT FRAMEWORK

2.1. Overall Architecture

The overall pipeline of the proposed solution is illustrated in Fig. 2. The framework includes three modular parts. First, the scale adapta-

tion module, composed of a tiny ResizeParamNet [21], a ConvNeXt-based pre/post-processor, and a bicubic interpolation filter. Second, the analysis/synthesis transform (g_a, g_s) of our design consists of a combination of patch merging/expanding layers and SwinT [24] blocks. The architectures of hyper-transforms (h_a, h_s) are similar to (g_a, g_s) with different stages and configurations. Finally, a Transformer-based slice transform under a ChARM design is used to estimate the distribution parameters of the quantized latent. These resulting discrete-valued data (\hat{y}, \hat{z}) are encoded into bit-streams with an arithmetic coder.

2.2. Scale Adaptation Module

Given a source image $x \in R^{H \times W \times C}$, we first determine an adaptive resize factor M estimated by the ResizeParamNet module, which consists of three stages of residual blocks (ResBlocks). Indeed, the estimated resize parameter M is used to create a sampling grid τ_M following the convention STNs, and used to adaptively down-scale x into $x_d \in R^{H' \times W' \times C}$ through the bicubic interpolation. The latter is then encoded and decoded with the proposed ICT. Finally, the decoded image $\hat{x}_d \in R^{H' \times W' \times C}$ is up-scaled to the original resolution $\hat{x} \in R^{H \times W \times C}$ using the same, initially estimated, resize parameter M . The parameterization of each layer is detailed in the ResizeParamNet and ResBlock diagrams of Fig. 3 (a) and (b), respectively. In addition, a learnable depth-wise pre/post-processor is placed before/after the bicubic sampler to mitigate the information loss introduced by down/up-scaling, allowing the retention of information. This neural pre/post-processing method consists of concatenation between the input and the output of three successive ConvNeXt [25] blocks. The ConvNeXt block diagram is also illustrated in Fig. 3 (c). For a better complexity-efficient design, we decided to skip the scale adaptation module where $M \cong 1$.

2.3. Transformer-based Analysis/Synthesis Transform

The analysis transform g_a contains four stages of patch merging layer and SwinT block to obtain a more compact low-dimensional latent representation y . In order to consciously and subtly balance the importance of feature compression through the end-to-end learning framework, we used two additional stages of patch merging layer and SwinT block in the hyper-analysis transform to produce an additional latent representation z . During training, both latents y and z are quantized using a rounding function to produce \hat{y} and \hat{z} , respectively. The quantized latent variables \hat{y} and \hat{z} are then entropy coded regarding an indexed entropy model for a location-scale family of random variables parameterized by the output of the ChARM, and a batched entropy model for continuous random variables, respectively, to obtain the bit-streams. Finally, quantized latents \hat{y} and \hat{z} feed the synthesis and hyper-synthesis transforms, respectively, to generate the reconstructed image. The decoder schemes are symmetric to those of the encoder, with patch-merging layers replaced by patch-expanding layers.

2.4. Transformer-based Slice Transform

Although there are strong correlations among different channels in latent space, the strongest correlations may come from the spatio-channel dependencies. Thus, to better parameterize the distribution of the quantized latents with a more accurate and flexible entropy model and without increasing the compression rate, we propose a Transformer-based slice transform inside the ChARM. Unlike previous works, ours considers spatio-channel latent correlations for entropy modeling in an auto-regressive manner. As a side effect,

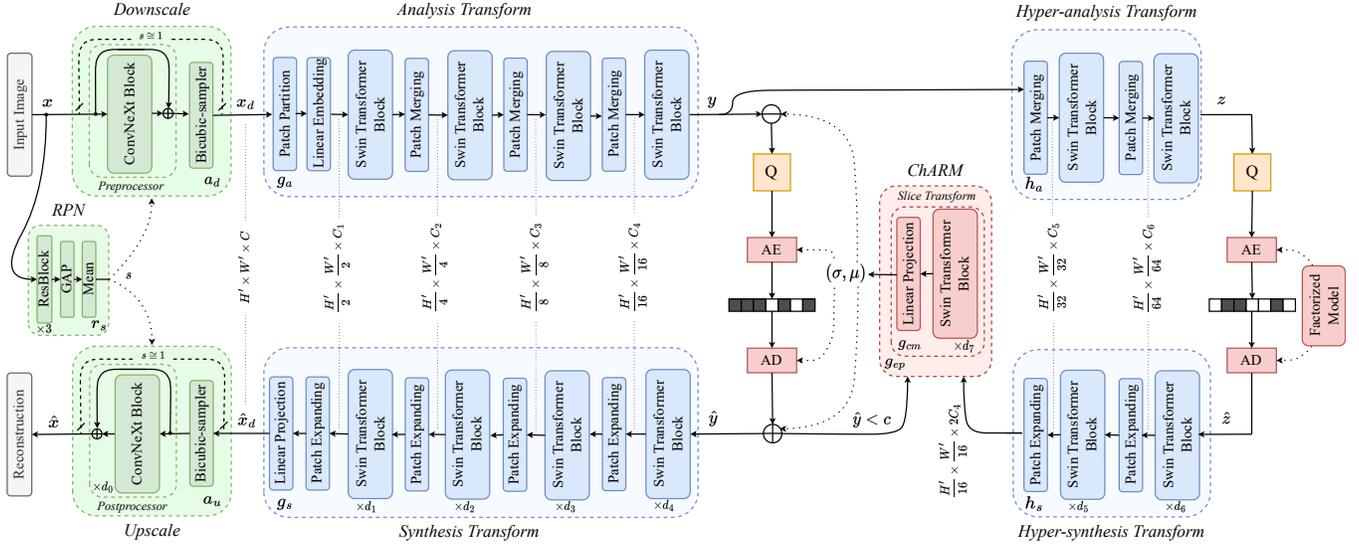


Fig. 2. Overall AICT framework. We illustrate the image compression diagram of our AICT with hyperprior and SwinT-based ChARM, and scale adaptation module. The ResizeParamNet and ConvNeXt block diagrams are illustrated in Fig. 3 (a) and (c).

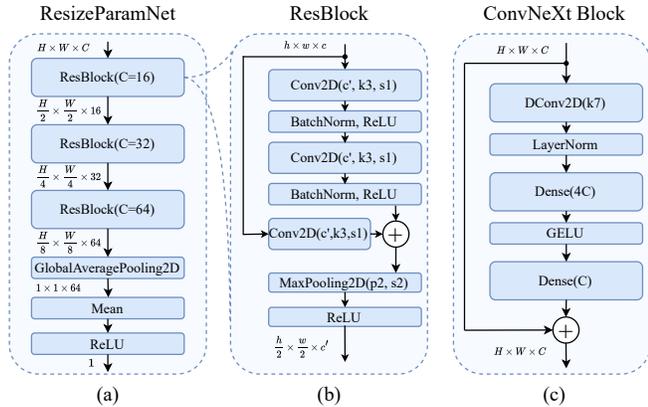


Fig. 3. Adaptation module block architectures.

it also leads to faster decoding speed. The slice transform consists of two successive SwinT blocks with an additional learnable linear projection layer, used to get a representative latent slices concatenation. This ChARM estimates the distribution $p_{\hat{y}}(\hat{y}|\hat{z})$ with both the mean and standard deviation of each latent slice, and incorporates an auto-regressive context model to condition the already-decoded latent slices and further reduce the spatial redundancy between adjacent pixels.

3. RESULTS AND ANALYSIS

3.1. Experimental Setup

Baselines.¹ We compare our solution with the state-of-art neural codec SwinT-ChARM proposed by Zhu *et al.* [14], and the Conv-ChARM proposed by Minnen *et al.* [13] and conventional codecs, including better portable graphics (BPG)(4:4:4), and the versatile

¹For a fair comparison, we only considered SwinT-ChARM [14] from the state-of-the-art models [14–19], due to the technical feasibility of models training and evaluation under the same conditions and in an adequate time.

video coding (VVC) official Test Model VTM-18.0 in All-Intra configuration.

Implementation details. We implemented all models in TensorFlow using tensorflow compression (TFC) library, and the experimental study was carried out on an RTX 5000 Ti GPU and an Intel(R) Xeon(R) W-2145 @ 3.70GHz CPU. All models were trained on the same CLIC20 training set with 2M iterations 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 the last 200k iterations, and $L = D + \lambda R$ is used as a loss function. L is a weighted combination of bitrate R and distortion D , with λ being the Lagrangian multiplier steering RD trade-off. Mean squared error (MSE) is used as the distortion metric in RGB color space. Each training batch contains eight random crops $\in R^{256 \times 256 \times 3}$ from the CLIC20 training set. To cover a wide range of rate and distortion points, for our proposed method and respective ablation models, we trained four models with $\lambda \in \{1000, 200, 20, 3\} \times 10^{-5}$. We evaluate the image codecs on four datasets [26], including Kodak, Tecnick, JPEG-AI, and the testing set of CLIC21. For a fair comparison, all images are cropped to the highest possible multiples of 256 to avoid padding for neural codecs.

3.2. Rate-Distortion Coding Performance

To demonstrate the compression efficiency of our proposed approach, we summarize, in Table 1, the BD-rate of our models and the baselines across four datasets compared to the VTM-18.0 as the anchor. On average, AICT is able to achieve 5.11% BD-rate reduction compared to VTM-18.0 and 3.93% relative gain from SwinT-ChARM. Also, we illustrate in Figure 4 a comparison of compression efficiency on Kodak dataset. 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 ICT and AICT achieve a good trade-off between BD-rate performance and decoding time.

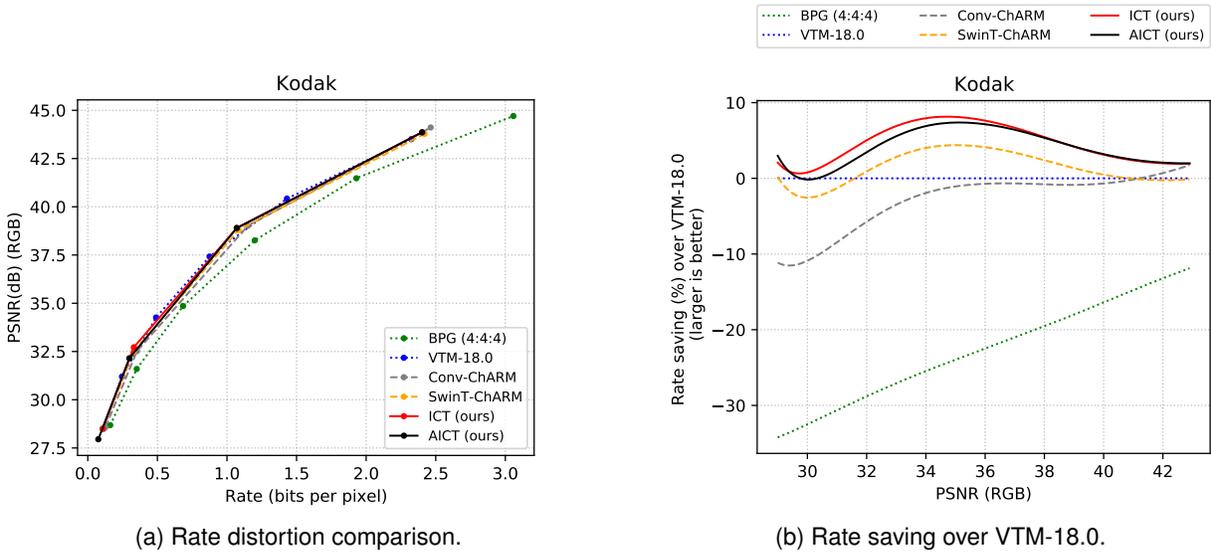


Fig. 4. Compression efficiency comparison on the Kodak dataset.

Table 1. BD-rate \downarrow performance of BPG (4:4:4), Conv-ChARM, SwinT-ChARM, ICT, and AICT compared to the VTM-18.0.

Image Codec	Kodak	Tecnick	JPEG-AI	CLIC21	Average
BPG444	22.28%	28.02%	28.37%	28.02%	26.67%
Conv-ChARM	2.58%	3.72%	9.66%	2.14%	4.53%
SwinT-ChARM	-1.92%	-2.50%	2.91%	-3.22%	-1.18%
ICT (ours)	-5.10%	-5.91%	-1.14%	-6.44%	-4.65%
AICT (ours)	-5.09%	-5.99%	-2.03%	-7.33%	-5.11%

3.3. Models Scaling Study

We evaluated the decoding complexity of the four considered image codecs by averaging decoding time across 7000 images encoded at 0.8 bpp. We present the image codecs complexity in Table 2, including decoding time on GPU and CPU, codec floating point operations per second (FLOPs), and codec total number of parameters. Compared to the neural baselines, ICT can achieve faster decoding speed on GPU but not on CPU, which proves the parallel processing ability to speed up compression on GPU and the well-engineered designs of both transform and entropy coding, highlighting an efficient and hardware-friendly compression model. This is potentially helpful for conducting high-quality real-time visual data streaming. Our AICT is on par with ICT in terms of the number of parameters, FLOPs, and latency, indicating that the scale adaptation module is not computationally heavy for real scenario applications.

3.4. Ablation Study

To investigate the impact of the proposed ICT and AICT, we conduct an ablation analysis according to the reported BD-rate results in Table 1. The compression performance increases from Conv-ChARM to SwinT-ChARM on the considered datasets due to the inter-layer feature propagation across non-overlapping windows (local information) and self-attention mechanism (local information) in the SwinT. With the proposed spatio-channel entropy model, ICT is able to achieve, on average, -3.47% BD-rate reduction com-

Table 2. Average decoding latency across 7000 images at 256×256 resolution, encoded at 0.8 bpp.

Image Codec	Latency(ms) \downarrow		MFLOPs \downarrow	#parameters (M) \downarrow
	GPU	CPU		
Conv-ChARM	133.8	359.8	126.1999	53.8769
SwinT-ChARM	91.8	430.7	63.2143	31.3299
ICT (ours)	80.1	477.0	74.7941	37.1324
AICT (ours)	88.3	493.3	74.9485	37.2304

pared to SwinT-ChARM. Therefore, introducing the Transformer-based slice transform leads to significant improvement compared to the ConvNet-based entropy model using only short-range dependencies. In addition, our spatio-channel entropy model is more helpful when combined with the Transformer-based transform coding. AICT performs better than ICT, indicating that the introduction of a scale adaptation module can further reduce spatial redundancies and alleviate coding artifacts, especially at low bitrate resulting in higher compression efficiency.

4. CONCLUSION

In this paper, we have proposed an up-and-coming neural codec AICT, achieving compelling RD performance while significantly reducing the latency, which is potentially helpful to conduct, with further optimizations, high-quality real-time visual data compression. We inherited the advantages of self-attention units from Transformers to effectively approximate both the mean and standard deviation for entropy modeling and combine global and local texture to capture correlations among spatially neighboring components for non-linear transform coding, achieving -4.65% BD-rate reduction over the VTM-18.0, by averaging over the benchmark datasets. Furthermore, we presented a lightweight scale adaptation module to enhance compression ability, especially at low bitrates, reaching on average -5.11% BD-rate reduction over the VTM-18.0.

5. REFERENCES

- [1] George Toderici, Sean M O'Malley, Sung Jin Hwang, Damien Vincent, David Minnen, Shumeet Baluja, Michele Covell, and Rahul Sukthankar, "Variable rate image compression with recurrent neural networks," *arXiv preprint arXiv:1511.06085*, 2015.
- [2] 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.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] Yin hao Zhu, Yang Yang, and Taco Cohen, "Transformer-based transform coding," in *International Conference on Learning Representations*, 2021.
- [15] 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.
- [16] 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.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al., "Spatial transformer networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [21] Li-Heng Chen, Christos G Bampis, Zhi Li, Lukáš Krásula, and Alan C Bovik, "Estimating the resize parameter in end-to-end learned image compression," *arXiv preprint arXiv:2204.12022*, 2022.
- [22] Hossein Talebi and Peyman Milanfar, "Learning to resize images for computer vision tasks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 497–506.
- [23] Chen Jin, Ryutarō Tanno, Thomy Mertzanidou, Eleftheria Panagiotaki, and Daniel C Alexander, "Learning to downsample for segmentation of ultra-high resolution images," *arXiv preprint arXiv:2109.11071*, 2021.
- [24] 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.
- [25] 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.
- [26] Benchmark datasets, "kodak testing set: <http://r0k.us/graphics>, technick testing set: <https://testimages.org/>, jpeg-ai testing set: https://jpegai.github.io/test_images/, and clic21 testing set: <http://compression.cc/tasks/>,".