

Transformer-based Image Compression

Ming Lu[†], Peiyao Guo[†], Huiqing Shi[‡], Chuntong Cao[‡], and Zhan Ma[†]

[†]Nanjing University, [‡]Jiangsu Longyuan Zhenhua Marine Engineering Co.

Abstract

A Transformer-based Image Compression (TIC) approach is developed which reuses the canonical variational autoencoder (VAE) architecture with paired main and hyper encoder-decoders. Both main and hyper encoders are comprised of a sequence of neural transformation units (NTUs) to analyse and aggregate important information for more compact representation of input image, while the decoders mirror the encoder-side operations to generate pixel-domain image reconstruction from the compressed bitstream. Each NTU is consist of a Swin Transformer Block (STB) and a convolutional layer (Conv) to best embed both long-range and short-range information; In the meantime, a casual attention module (CAM) is devised for adaptive context modeling of latent features to utilize both hyper and autoregressive priors. The TIC rivals with state-of-the-art approaches including deep convolutional neural networks (CNNs) based learnt image coding (LIC) methods and handcrafted rules-based intra profile of recently-approved Versatile Video Coding (VVC) standard, and requires much less model parameters, e.g., up to 45% reduction to leading-performance LIC.

1 Introduction

High-efficiency image compression plays a vital role for effective Internet service, such as the online advertisements, professional photography sharing, etc. Though traditional rules-based image coding standards, e.g., JPEG [1], JPEG2000 [2], Video Coding Intra Profile, etc, are widely deployed for decades, a better image coding approach, i.e., less bandwidth consumption but better reconstruction quality, is a constant desire for better service.

Recent years, a variety of deep learning based image compression approaches [3, 4, 5, 6, 7] have emerged with noticeable coding efficiency improvement, even offering better performance than the VVC Intra [8] quantitatively and qualitatively [6, 7]. This attractive potentiality encourages the pursuit of next-generation image coding techniques from both academia researchers and industrial leaders. Thus, international standardization groups, such as the well-known ISO/IEC JPEG, have officially called for the technical evaluation and standardization of deep learning based image compression.

1.1 Background and Motivation

As seen, existing LICs generally use CNNs in VAE framework where a number of convolutional layers (with resolution scaling) are stacked to exploit spatial correlation in local neighborhood to derive compact latent features at bottleneck for entropy coding. Applying CNNs for discriminative feature extraction was originally inspired by the discoveries of Hubel and Wiesel about the cats' visual cortex in 1960s [9].

However, the CNNs do have limitations [10]. First, convolutional filter can only characterize short-range spatial correlation within the receptive field; Then, offline trained CNN model uses fixed parameters, making it generally incapable of dealing with images having very different content distribution; Third but not the last, CNN computation is computational intensive due to element-by-element processing.

To overcome the drawbacks of native CNNs, a number of principled rules are then developed. For example a variety of attention mechanisms, e.g., nonlocal, spatial- and channel-wise methods, are manually integrated with convolutions to capture long-range correlations for more compact latent feature with better compression performance [7, 6]. However, these handcrafted attention methods often require excessive model parameters which makes them impractical for real-life applications. To deal with images with different content distributions, multi-model optimization helps the LIC encoder to select an optimal one with better rate-distortion measurements [11], which apparently requires us to cache multiple models.

Recalling that the human visual system (HVS) usually scans the natural environment during a saccade to fixate on different regions for saliency extraction and scene understanding, it helps us to effectively capture long-range correlations for better information embedding. Such biological visual fixation with attentive window adaptation is well simulated by the shifted window method suggested in award-winning Swin Transformer[12]. And, having the successful deployment of convolutions in existing LICs to embed local neighborhood information, this work suggests to combine the convolutional layer (for short-range information embedding) and Swin Transformer-based attention block (for long-range information aggregation), by which we wish to produce more compact latent features for better coding efficiency.

1.2 Our Approach and Contribution

This paper proposes the Transformer-based Image Compression (TIC) method. The TIC applies the same VAE architecture as used in existing LICs. For better information embedding, we suggest the neural transformation unit (NTU) as the basic module which is comprised of a Swin Transformer block (STB) and a convolutional layer (Conv). Similarly, we use the resolution scaling in convolutional layer (with pre-defined strides) to characterize and embed spatial information. To efficiently encode the latent feature at the bottleneck, both hyper priors and autoregressive neighbors are utilized as in [5, 6, 7] through a causal attention module (CAM) that combines the causal self-attention and the multi-layer perceptron (MLP) to exploit closely-related priors for context modeling.

Experimental results show that the proposed TIC shows competitive compression efficiency, and only requires about a half of model parameters to the state-of-the-art LIC method - Cheng *et al.* [6]. We then extend the TIC by placing more STBs in layers closer to the bottleneck at main coder, denoted as the TIC+, by which we can even outperform the VVC Intra with a slight increase of model parameters.

2 Related Work

Learnt Image Compression. The core problem of image compression is about the compact representation of pixel blocks, for which over decades, the transforms, such as the Discrete Cosine Transform, Wavelet, etc, and context adaptive entropy codes have been extensively studied to fulfill the purpose.

Back to 2016, Ballé *et al.* [3] showed that stacked convolutional layers can replace the traditional transforms to form an end-to-end trainable image compression method with better efficiency than the JPEG. Then, by adopting the hyper prior [4] and autoregressive neighbors [5] for entropy context modeling, the image compression efficiency was further improved and was competitive to the Intra Profile of the High-Efficiency Video Coding (HEVC). Cheng *et al.* [6] later introduced a Gaussian mixture model for better approximating the distribution of latent features, with which the comparable performance to the VVC Intra was reported. In addition to these methods mainly utilizing the CNNs to analyze and aggregate information locally, our early exploration on neural image coding (NIC) in [7] applied a nonlocal attention to the intermediate features generated by each convolutional layer to select attentive features for more compact representation. However, the nonlocal computation is expensive since it typically requires a large amount of space to host a correlation matrix with size of $HW \times HW$ for a naive implementation. Note that H and W are the height and width for input feature map.

As seen, handcrafted attention mechanisms have already promised the potentials for improving the coding efficiency [6, 7], in which they are mostly utilized to guide the encoding of importance area. Whereas, this work suggests to combine the convolution and attention to jointly consider the short-range and long-range correlation effectively.

Vision Transformer. Since the AlexNet, the CNNs have improved the performance of numerous vision tasks remarkably, by relying on the powerful capacity of stacked convolutions in characterizing the input data. Even though, the locality of the convolutional computation limits its performance by only aggregating the information within the limited receptive field. Recently, emerging vision transformer networks [13, 14, 15] have demonstrated very encouraging improvement to those methods only using CNNs for various high-level vision tasks. One potential reason is that the Transformer structure can well capture the long-range correlation for better information embedding, by dividing the input images into tokenized patches and treating them as sequential words as in natural language processing for computation [16].

Later, a number of improvements have been made, including the position encoding, token pooling, feed forward networks and layer normalization, to further the performance. Particularly, the emergence of the award-winning Swin Transformer [12] has revealed that the Transformer can be also extended to low-level vision tasks [17, 10] for better reconstructed quality. The Swin Transformer applies the window-based attention and relative position encoding, making it capable of processing arbitrary-size high resolution images. Applying the computation on non-overlapped windows can significantly reduce the computational complexity when compared with the normal convolutions that traverses all elements in a frame. In addition, the shifted window scheme performs the nonlocal processing at the patch-sized level to exploit the

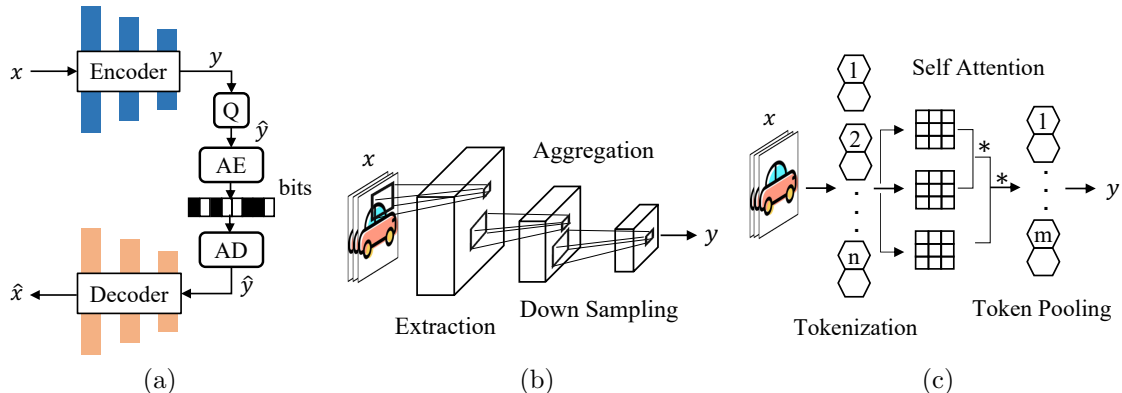


Figure 1: (a) The common paradigm of end-to-end learning-based image compression. Q, AE and AD represent the quantization, arithmetic encoding and decoding respectively. (b) The encoding flow of CNN-based image compression. (c) The encoding folow of Transformer-based image compression.

long-range dependency for better information aggregation.

This work attempts to migrate the Swin Transformer into the image compression pipeline for better performance.

3 Method

3.1 Local Convolution and NonLocal Transformer for Image Compression

Figure 1a illustrates the common paradigm of end-to-end learning-based image compression methods. The encoder transforms the input x into the latent features y , and then quantize y to discrete symbols \hat{y} that are entropy coded into bitstreams using predefined distribution models; The decoding part mirrors the encoding steps by parsing the compressed bitstream to reconstruct pixel blocks to form decoded \hat{x} . The optimization objective is to minimize the rate-distortion cost through an end-to-end learning means:

$$L = R(\hat{y}) + \lambda D(x, \hat{x}), \quad (1)$$

where R is compressed bit rate of \hat{y} and the distortion D measures the mean square error (MSE) between the ground truth x and restored output \hat{x} . We adapt λ for rate-distortion trade-off at various bit rates.

Figure 1b details the encoding flow of existing LICs that mainly use CNNs for image coding. It uses stacked convolutional layers to analyze and aggregate features for compact representation of the input at the bottleneck. Often times, convolution is coupled with predefined strides for resolution scaling and spatial neighborhood information embedding. As seen, the receptive field can be enlarged by the resolution scaling to exploit spatial correlation from more local neighbors. At the bottleneck, quantized latent features are then entropy-coded through context models conditioned on autoregressive neighbors and hyper priors. To jointly utilizing the autoregressive

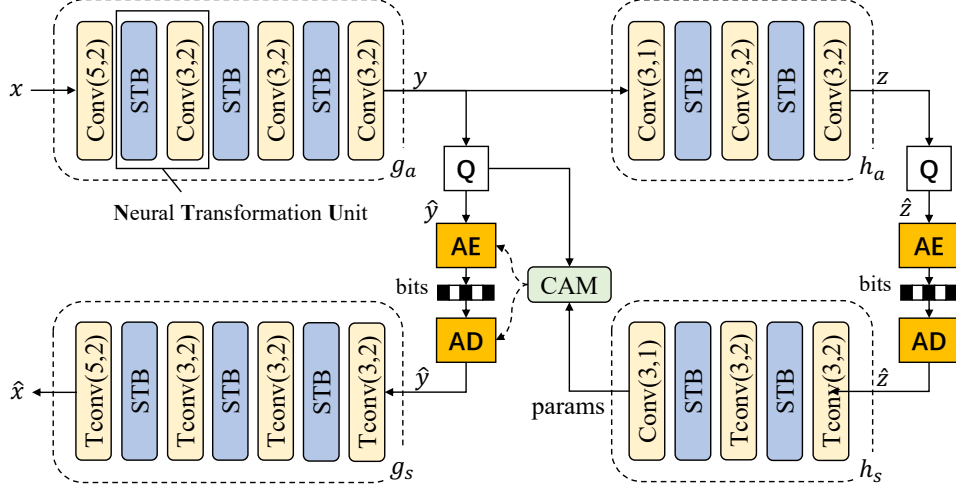


Figure 2: **TIC**. The proposed TIC stacks the Neural Transformation Unit (NTUs) in both main and hyper encoder-decoders of a VAE framework. Each NTU is comprised of a STB and a Conv layer for which we wish to capture and embed long-range and short-range information. Conv(k,s) and Tconv(k,s) are convolutional and transpose convolutional layers with kernel size $k \times k$ and stride size s . $k = 3$ and $s = 2$ are exemplified in this work. The causal attention module (CAM) is applied to aggregate information from autoregressive neighbors and hyper priors for context modeling which is different from existing masked CNN based information fusion mechanism.

neighbors and hyper priors for better context modeling, masked CNN [18] is typically used to fuse information locally.

Additionally, Figure 1c shows the encoding flow by using the Transformer. Input image is first tokenized using a convolutional layer to produce fixed-size tokens (i.e., feature patches after sequential projection); Then self attention layer is applied to derive spatial coefficients for upcoming token pooling by which we can intentionally remove less important tokens.

This work attempts to leverage the advantages of both local convolution and Transformer-based nonlocal attention for better information embedding.

3.2 TIC: Exploring the Combination of Convolution & Transformer

Figure 2 details the diagram of the proposed TIC. We follow the canonical VAE architecture [4] to construct main and hyper encoder-decoder pairs. For the main encoder g_a , input image is first convoluted prior to being fed into succeeding three NTUs. The hyper encoder h_a shares similar architecture but with two NTUs. The main decoder g_s and hyper decoder h_s reverse the processing steps of g_a and h_a respectively. Each NTU applies a STB and a Conv layer to step-wisely analyze and embed respective long-range and short-range information.

As in STB, the first feature embedding (FE) layer projects input features at a size of $H \times W \times C$ to a dimension of $HW \times C$, the following Swin Transformer

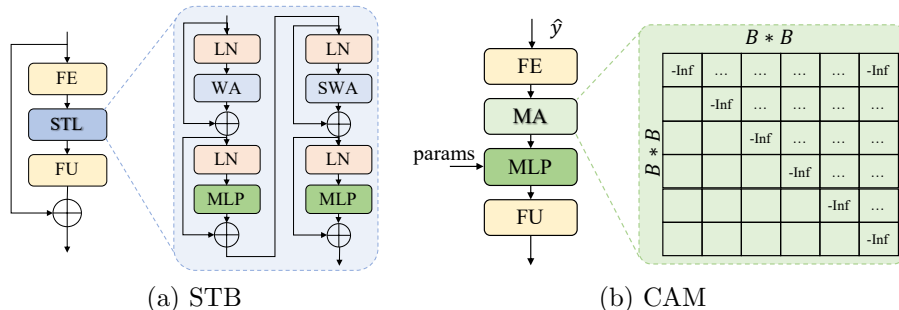


Figure 3: **TIC Modules.** (a) Swin Transformer Block (STB); (b) Causal Attention Module (CAM) for context modeling.

layer (STL) which consists of the layer normalization (LN), window attention (WA), shifted window attention (SWA) and the MLP layer calculates the window-based self-attention, and finally a feature unembedding (FU) layer remaps attention weighted features back to original size of $H \times W \times C$. Skip connection is used for better aggregation. Note that no patch division and fully-connected layers for tokenization in our work are beneficial to early visual processing and stable training [19, 20].

As in [7], resolution scaling is integrated in the Conv layer with stride 2, by which we further aggregate local blocks or tokens in a spatial neighborhood. Currently, we simply use the linear GDN and Leaky ReLU for activation after the Conv in main and hyper coders. Other simple activation functions such as the ReLU can be used as well.

Different from existing masked CNN based context modeling, we propose a causal attention module (CAM) in Fig. 3b to select attentive neighbors from autoregressive priors where the CAM unfolds the quantized features into $B \times B$ patches and calculate relations among these patches with masked attention (MA) to guarantee the causality. Experiments suggest $B = 5$ for well balancing the performance and complexity. The succeeding MLP layer fuse attention weighted autoregressive neighbors and hyper priors from h_s for final context prediction.

We later extend the TIC by placing more STBs at layers closer to the bottleneck, i.e., having 1 STB, 2 STBs, and 3 STBs respectively in three NTUs of the main encoder-decoder, and keeping the same 1 STB in each NTU in hyper coder, named as the TIC+.

4 Results

Experimental Setup. The Flickr 2W [21] is used as the training dataset. We randomly crop images into fixed patches at a size of $256 \times 256 \times 3$. Adam is used as the optimizer and the batch size is set to 8 for training. All training threads run on a single Titan RTX GPU for 400 epochs with the learning rate of 10^{-4} . We now train 8 models in total to match different bit rates (or quality levels) that can be adapted by selecting a λ is $\{0.0018, 0.0035, 0.0067, 0.013, 0.025, 0.0483, 0.0932, 0.18\}$. The

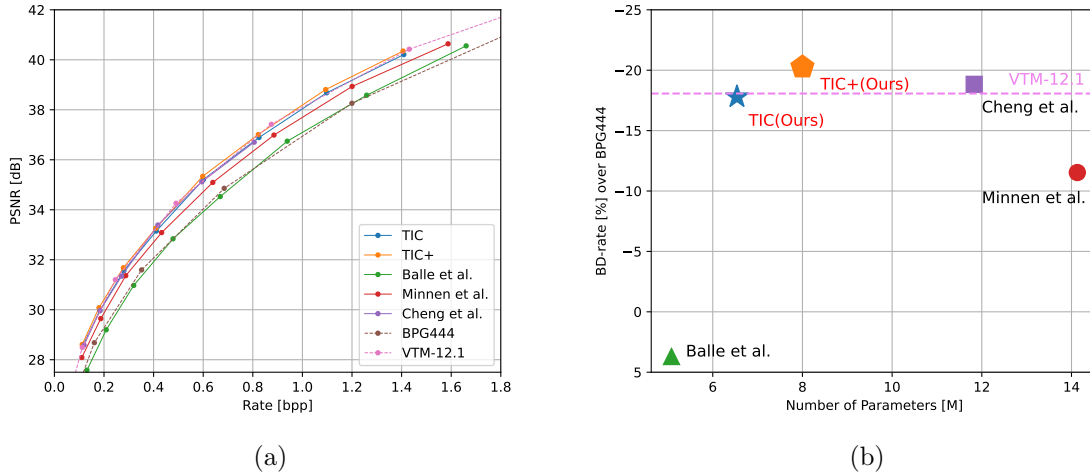


Figure 4: **Quantitative Evaluation.** (a) R-D Performance of the TIC and TIC+ against the BPG444 (HEVC Intra), VTM 12.1 (VVC Intra) and other LICs with leading performance [4, 5, 6]; (b) Performance gains (BD-Rate) against the HEVC Intra using BPG444 mode, and model parameters (in Mbytes). Upper left is better. Note that performance is averaged using all test images in Kodark dataset.

MSE is used for distortion measurements.

The proposed TIC is implemented on top of an open-source CompressAI PyTorch library [22], by which we can easily share our models and materials for reproducible research. We evaluate our model on the Kodak dataset, having the peak signal-to-noise ratio (PSNR) to quantify the image quality and the bits per pixel (bpp) to measure the bit rate.

In STB, the window sizes are set to 8×8 and 4×4 in respective main and hyper encoder-decoders and the numbers of heads are 4, 8, 16 for three STBs in main encoder-decoders and 16 for hyper STBs. We adopt 128 channels for the first 4 models corresponding to low bit rates scenarios, and 192 channels for the rest 4 models to cover high bit rates.

Quantitative and Qualitative Evaluations. Quantitative performance illustrated either using rate-distortion (R-D) curves in Figure 4a or using BD-Rate gains (over the HEVC Intra) in Figure 4b reports the competitive efficiency of the proposed TIC to the state-of-the-art Cheng *et al.* [6], and VVC Intra (VTM 12.1), but the TIC only requires a half of model parameters to Cheng *et al.* [6] which is more preferred for real-life application.

And, the TIC+ even surpasses the VVC Intra by 2.6% BD-Rate improvement with a slight increase of model parameters. Even though, the TIC+ still consumes much less parameters than the Cheng *et al.* [6].

Figure 5 visualize the reconstructions and closeups generated by the TIC, HEVC Intra (BPG), and VVC Intra (VTM), which agrees with the objective improvement

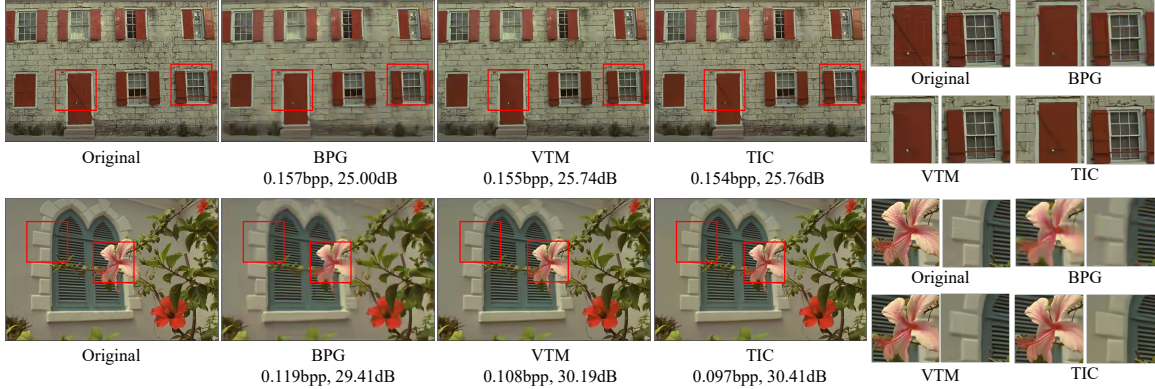


Figure 5: **Qualitative Visualization.** Reconstructions and close-ups of the TIC, BPG (HEVC Intra) & VTM (VVC Intra). Corresponding bpp and PSNR are marked.

of proposed TIC by subjective illustration with more sharp textures and less noise.

All of these studies have revealed the efficiency by combing the convolution and attention-based Transformer. By placing more STBs in TIC+, we observe more coding gains. This would be an interesting topic in future to study how many STBs are sufficient for achieving the optimal coding efficiency.

As aforementioned, the core issue of image compression is about the generation of compact representation of input image. We then use the partial derivative $\frac{\partial g_a(m,n)}{\partial I(i,j)}$ to visualize the compactness of input pixel $I(i,j)$ to the latent feature on position (m,n) generated by the main encoder g_a . Figure 6 exemplifies the scenario when we take the center of g_a to illustrate the contribution from all pixels of input image. Apparently, our TIC offers much compact illustration than Cheng *et al.*'s model that uses the CNN-based method, which further evidence that the combination of Transformer-based attention and convolution can better embed spatial information with more compact representation, and thus can lead to better coding efficiency as reported.

5 Conclusion

This paper reports a state-of-the-art image compression method. It combines the Swin Transformer and the convolutional layer as the basic unit to analyze and aggregate short-range and long-range information for more compact representation of input image. Experimental results reveal the leading performance when compared with existing learning-based approaches, and recently-emerged VVC Intra. Besides the encouraging coding efficiency, the proposed method consumes much less model parameters to the existing learning-based approaches, making the solution attractive to practical applications. An interesting topic for future study is to further improve the coding efficiency along the direction by stacking the convolutions and Transformers. More test results and models will be updated regularly at <https://njuvision.github.io/TIC>.

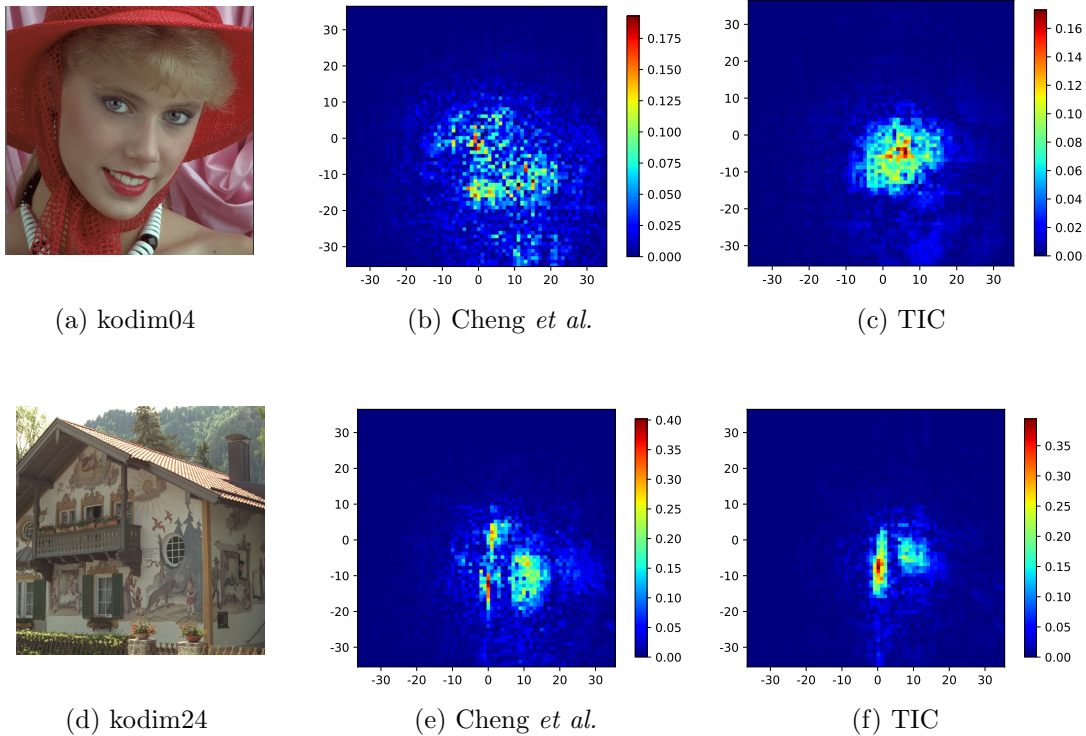


Figure 6: **Compactness Visualization of Latent Features.** Close-ups of the gradient maps averaged over all channels for different test images.

References

- [1] Gregory K Wallace, “The jpeg still picture compression standard,” *IEEE transactions on consumer electronics*, vol. 38, no. 1, pp. xviii–xxxiv, 1992.
- [2] Majid Rabbani, “Jpeg2000: Image compression fundamentals, standards and practice,” *Journal of Electronic Imaging*, vol. 11, no. 2, pp. 286, 2002.
- [3] Johannes Ballé, Valero Laparra, and Eero P Simoncelli, “End-to-end optimized image compression,” in *International Conference on Learning Representations*, 2017.
- [4] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston, “Variational image compression with a scale hyperprior,” in *International Conference on Learning Representations*, 2018.
- [5] David Minnen, Johannes Ballé, and George Toderici, “Joint autoregressive and hierarchical priors for learned image compression,” in *Advances in Neural Information Processing Systems*, 2018, pp. 10794–10803.
- [6] 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.
- [7] 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.

- [8] 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.
- [9] David H Hubel and Torsten N Wiesel, “Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex,” *The Journal of physiology*, vol. 160, no. 1, pp. 106–154, 1962.
- [10] Jingyun Liang, Jie Zhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte, “Swinir: Image restoration using swin transformer,” in *IEEE International Conference on Computer Vision Workshops*, 2021.
- [11] Guo Lu, Xiaoyun Zhang, Wanli Ouyang, Li Chen, Zhiyong Gao, and Dong Xu, “An end-to-end learning framework for video compression,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [12] 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,” *arXiv:2103.14030*, 2021.
- [13] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko, “End-to-end object detection with transformers,” in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [15] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng, “Deepvit: Towards deeper vision transformer,” *arXiv:2103.11886*, 2021.
- [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [17] Zhendong Wang, Xiaodong Cun, Jianmin Bao, and Jianzhuang Liu, “Uformer: A general u-shaped transformer for image restoration,” *arXiv:2106.03106*, 2021.
- [18] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma, “Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications,” *arXiv preprint arXiv:1701.05517*, 2017.
- [19] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [20] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick, “Early convolutions help transformers see better,” *arXiv preprint arXiv:2106.14881*, 2021.
- [21] Jiaheng Liu, Guo Lu, Zhihao Hu, and Dong Xu, “A unified end-to-end framework for efficient deep image compression,” *arXiv preprint arXiv:2002.03370*, 2020.
- [22] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja, “Compressai: a pytorch library and evaluation platform for end-to-end compression research,” *arXiv preprint arXiv:2011.03029*, 2020.