

Adaptive Learned Image Compression with Graph Neural Networks

Supplementary Material

001 1. Advanced CNNs and Transformers

002 We provide additional quantitative comparisons
003 against advanced CNN- and Transformer-based oper-
004 ators. Specifically, we replace our GNN module
005 with either deformable/dynamic convolutions [3, 4] or
006 Swin/Neighborhood Attention (NAT) [7, 13]. We keep the
007 entropy model and FFN unchanged, and only swap the
008 graph operator with these alternative operators. All models
009 are trained with $\lambda = 0.05$, yielding the loss curves in Fig. 1
010 and the R-D points in Fig. 2.

011 At $\lambda = 0.05$, GLIC reduces bitrate by more than 9%
012 compared with deformable/dynamic convolutional variants,
013 and by more than 5% compared with Swin/NAT variants,
confirming the advantage of our GNN-based design.

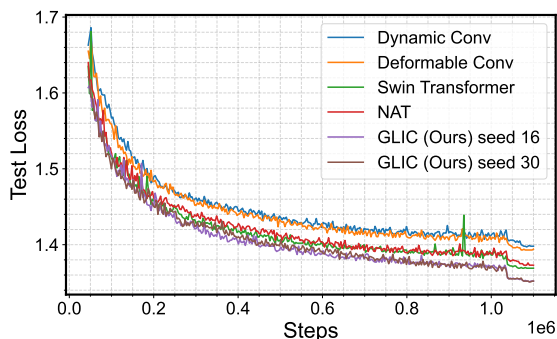


Figure 1. Test Loss on Kodak vs. Training Step

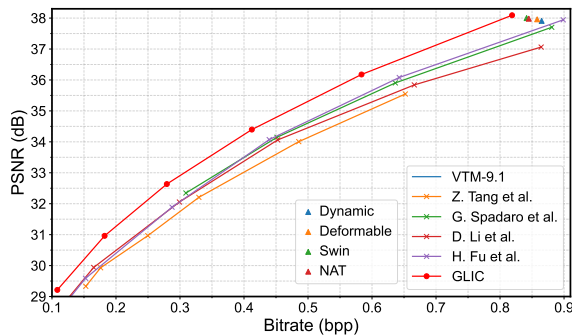


Figure 2. PSNR-based RD curves on Kodak dataset.

015 2. Training Stability

016 Our GLIC model can be trained stably under different ran-
017 dom seeds. The test loss curves of GLIC with multi-
018 ple seeds (Fig. 1) consistently indicate stable optimization.
019 Compared with other operators and baselines, GLIC ex-
020 hibits smooth convergence and consistently outperforms the
021 alternatives when using identical entropy models.

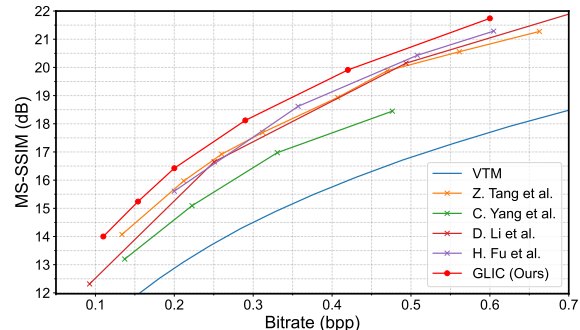


Figure 3. MS-SSIM-based RD curves on Kodak dataset.

022 3. Comparisons with GNN-based LIC models

023 We provide detailed PSNR and MS-SSIM R-D curves in
024 Fig. 2 and Fig. 3. We compare against three graph-based
025 LIC models, Z. Tang *et al.* [16], G. Spadaro *et al.* [15], and
026 C. Yang *et al.* [17], as well as two deformable CNN-based
027 LIC models, D. Li [9] and H. Fu *et al.* [6]. Our GLIC model
028 significantly outperforms all of them, demonstrating the ef-
029 fectiveness of our GNN design.

030 Our GLIC differs fundamentally from existing methods
031 (Tab. 1). These differences constitute the key innovations of
032 our approach:

033 1. Flexible Receptive Fields via Dual-scale Sampling.

034 Fully connected or global KNN graphs in [16, 17] provide
035 global context but incur $\mathcal{O}(N^2)$ complexity, while the
036 linear-complexity window graph in [15] restricts the recep-
037 tive field to a local patch. Our dual-scale graph combines
038 both types of graphs, achieving a near-global receptive field
039 with linear complexity.

040 2. Complexity-aware Neighbor Quota Assignment.

041 Instead of fixing the neighbor count to a constant K , we adap-
042 tively allocate node degrees according to the RMS gradient
043 magnitude, assigning denser connections to visually complex
044 regions.

045 4. Details of ERF Computation

046 The effective receptive field (ERF) [14] characterizes how
047 strongly each input pixel *actually* affects a designated out-
048 put unit, complementing the theoretical receptive field im-
049 plied by the architecture. Given an input image $\mathbf{x}_0 \in \mathbb{R}^{3 \times H \times W}$
050 and a network output (or intermediate feature) \mathbf{y} , we denote by
051 $x_{0,i,j}$ the input pixel at spatial location (i, j) , and by
052 $y_{c,h,w}$ the feature at channel c and spatial location
053 (h, w) .

054 Following our implementation, we compute the ERF on
055 the analysis transform $g_a(\cdot)$. For an input image, we first

Table 1. Comparison with other GNN-based models. N : total pixels, M : pixels per window [2], K : the GNN neighborhood size, n : sample size for our graphs. Complexity: *Construction + Aggregation* Complexity.

	[16]	[17]	[15]	GLIC (ours)
Construction	Fully-connected graph	Global K -NN graph	Local window K -NN	Dual graph
Node Degree	Fixed $N-1$ (full)	Fixed K	Fixed K	Adaptive, average n
Aggregation	Graph Attention	ECC	Graph Attention	Edge-conditioned Aggregation
Receptive field	Full	Full	Bounded by Window	Approximately Full
Complexity	$O(1) + O(N^2)$	$O(N^2) + O(NK)$	$O(NM) + O(NK)$	$O(Nn) + O(Nn)$

056 obtain the latent feature

$$057 \quad \mathbf{y} = g_a(\mathbf{x}_0), \quad (1)$$

058 and select its spatial center (h_c, w_c) . To measure the influ-
059 ence on the center latent unit, we form a scalar response by
060 summing all channels at that position:

$$061 \quad s = \sum_c y_{c, h_c, w_c}. \quad (2)$$

062 The ERF map is then defined as the magnitude of the gradi-
063 ent of s with respect to the input:

$$064 \quad \text{ERF}(i, j) = \left| \frac{\partial s}{\partial x_{0, i, j}} \right|. \quad (3)$$

065 Intuitively, a larger $\text{ERF}(i, j)$ indicates that the input pixel
066 $x_{0, i, j}$ contributes more strongly to the center latent re-
067 sponse.

068 In practice, we back-propagate from s to the input and
069 post-process the resulting gradients: (i) clamp large values
070 by an upper bound τ (we use $\tau=0.2$); (ii) apply ReLU to
071 retain only positive contributions; and (iii) aggregate over
072 channels to obtain a 2D spatial map:

$$073 \quad \tilde{\mathbf{E}}(i, j) = \sum_c \text{ReLU} \left(\min \left(\frac{\partial s}{\partial x_{0, c, i, j}}, \tau \right) \right). \quad (4)$$

074 5. VTM Setting and Commands

075 We follow [5] to ensure a widely used baseline. Specif-
076 ically, we directly adopt the R-D results reported in the
077 CompressAI repository [1], which are obtained using VTM-
078 9.1. For our VTM-9.1 evaluations, we use the QP set
079 $\{22, 27, 32, 37, 42, 47\}$, with the following command:

```
080 VTM/bin/EncoderAppStatic -i [input.yuv] \  
081 -c VTM/cfg/encoder_intra_vtm.cfg \  
082 -o [output.yuv] -b [output.bin] \  
083 -wdt [width] -hgt [height] -q [QP] \  
084 --InputBitDepth=8 -fr 1 -f 1 \  
085 --InputChromaFormat=444
```

Table 2. BD-rate (%) comparison between the GNN-based SC-
CTX entropy model (GNN-SCCTX) and the convolution-based
SCCTX entropy model in GLIC. Lower is better.

Method	Kodak(%)	Tecnick(%)	CLIC(%)	Dec.Lat(s)
GNN-SCCTX	-19.08	-22.86	-19.78	0.506
Conv-SCCTX (Ours)	-19.29	-21.69	-18.71	0.395

086 6. GNNs in the Entropy Model

087 Our experiments so far demonstrate that the proposed GNN
088 is highly effective for nonlinear transform networks. In the
089 main GLIC model, we follow [8] and adopt convolutional
090 networks and MLPs for spatial-channel entropy modeling.
091 Here, we further explore the potential of GNNs for improv-
092 ing the entropy model itself.

093 As shown in Tab. 2, we construct a variant (GNN-
094 SCCTX) that uses our GNNs to model the spatial and chan-
095 nel context in the SCCTX entropy model, and compare it
096 with the convolution-based SCCTX (Conv-SCCTX). We
097 observe that the GNN-SCCTX variant achieves better R-D
098 performance on the high-resolution datasets CLIC and Tec-
099 nick by more than 1.2%. However, this comes at the cost
100 of a 28% increase in decoding latency. We consider this
101 level of BD-rate gain insufficient to justify such a substan-
102 tial slowdown in decoding.

103 Nevertheless, these results indicate that the entropy
104 model can be further enhanced by GNNs or other advanced
105 architectures, suggesting that the overall performance of
106 GLIC could be further boosted when combined with a
107 stronger entropy model.

108 7. Ablations on GNN Settings

109 We further ablate the GNN design by varying the *aver-*
110 *age connection number* and the *candidate sampling size*, as
111 summarized in Tab. 3. Increasing these two hyperparam-
112 eters to “64 / 16×16” substantially improves the R-D per-
113 formance, while only mildly increasing the decoding latency.

114 We summarize the observation as follows: **(1) Sparsity.**
115 Increasing the candidate set size does not yield unbounded
116 RD gains. We find that 16×16 sampling, though sparse,
117 is already near saturation. A denser setting (e.g., 32×32)

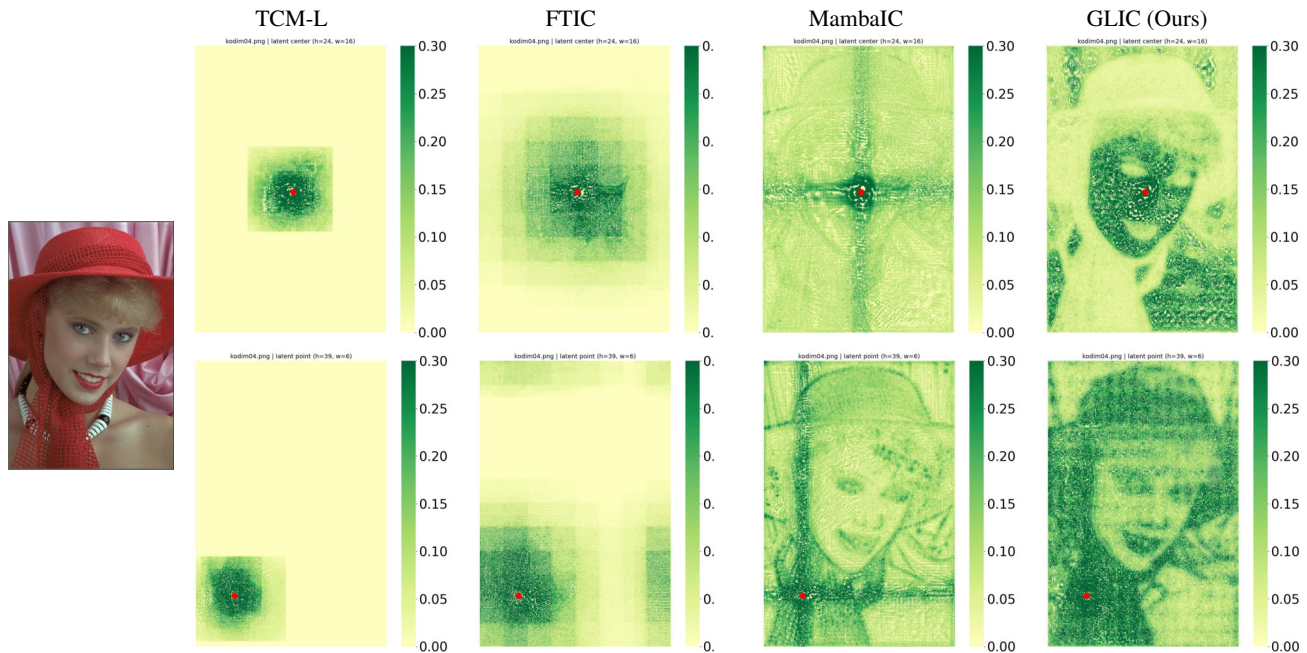


Figure 4. Additional single-image ERF visualizations (set 1).

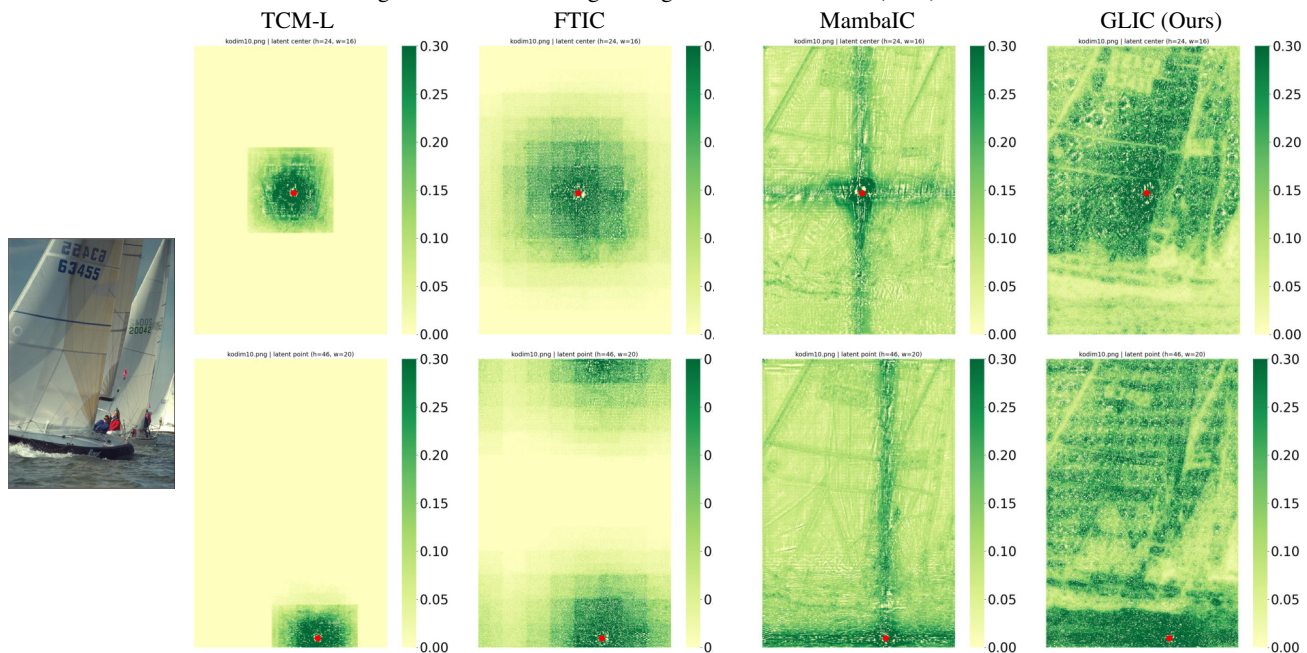


Figure 5. Additional single-image ERF visualizations (set 2).

118 yields no significant gain while incurring much higher latency. (2) **Budget.** Varying the connection budget (average
 119 degree) shows that 64 neighbors are sufficient. Expanding
 120 to a fully connected graph (256 connections) does not further
 121 improve RD performance as it may result in unnecessary
 122 interactions for redundancy elimination.
 123

8. Clarification on Table 3 in the Main Paper

In Table 3 (“Ablations on dual-scale graph flexibility”) of
 the main paper, the entries “w/o Global (Local only)” and
 “w/o Local (Global only)” indicate that one of the two
 graphs is replaced by the other one. In other words, each
 variant replaces the dual-scale design with a single-scale
 design, but does not change the network depth.

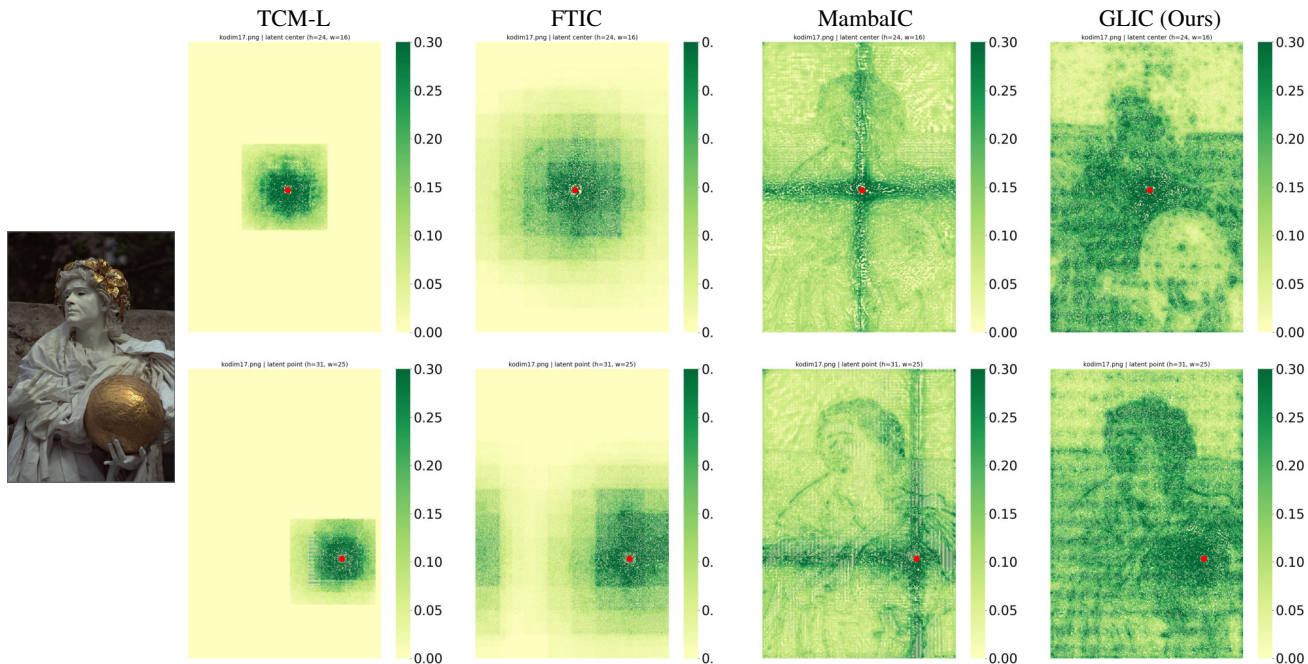


Figure 6. Additional single-image ERF visualizations (set 3).

Table 3. Ablation on average connection number and candidate sampling size. BD-rate (%) is computed on Kodak; lower is better.

Average Connection / Candidate Size	Kodak(%)	Dec.Lat(s)
32 / 8×8	-15.33	0.298
64 / 16×16 (Ours)	-19.29	0.395
256 / 16×16	-19.17	0.402
512 / 32×32	-19.97	0.891

Table 4. Comparison between standard and accelerated training under different λ values on Kodak (lower test loss is better).

	$\lambda = 0.05$	$\lambda = 0.0017$
Standard Training	1.344	0.258
Accelerated Training	1.343	0.258

131 9. Ablations on Training Strategy

132 We also ablate the accelerated training strategy proposed in
 133 [11] at two operating points, $\lambda = 0.05$ and $\lambda = 0.0017$. We
 134 compare the best test losses of GLIC obtained with this ac-
 135 celerated strategy against those obtained with the standard
 136 2M-step training schedule used in [8, 10, 12]. The results
 137 on Kodak are reported in Tab. 4.

138 The test losses under the two training regimes are almost
 139 identical, and the remaining gaps are very small. Therefore,
 140 we conclude that adopting the accelerated training strategy
 141 does not noticeably affect the performance of our GLIC
 142 model.

143 10. More Adaptivity Visualizations

144 We provide additional single-image ERF visualizations to
 145 further illustrate the adaptivity of our method. As shown in
 146 Fig. 4–6, GLIC produces ERFs that are both more global
 147 and more content-adaptive than those of prior approaches.
 148 The activated ERF patterns vary noticeably across different
 149 images and also across different target positions within the
 150 same image.

151 In contrast, previous Transformer- and Mamba-based
 152 methods, including TCM [12], FTIC [10], and Mam-
 153 baIC [18], tend to yield highly similar ERF structures across
 154 images and target locations. Their high-response regions
 155 are largely isotropic and appear content-agnostic rather than
 156 truly content-dependent.

157 Our method, however, exhibits clear content-driven be-
 158 havior: for example, in Fig. 4, it strongly attends to the face,
 159 the red hat and the scarf, respectively in the two examples;
 160 in Fig. 6, it focuses on the white sculpture and the golden
 161 sphere, respectively and separately for the two examples.
 162 Importantly, regions with large ERF values align with areas
 163 that exhibit substantial visual redundancy consistent with
 164 human perception and semantic structure, further validating
 165 the effectiveness of our approach.

166 11. More Visualization Comparisons

167 We also provide more subjective visual comparisons to
 168 demonstrate the effectiveness of our GLIC model. We com-
 169 pare GLIC with the traditional codec VTM-9.1 [2] and the
 170 LIC model TCM [12]. A detailed inspection of image

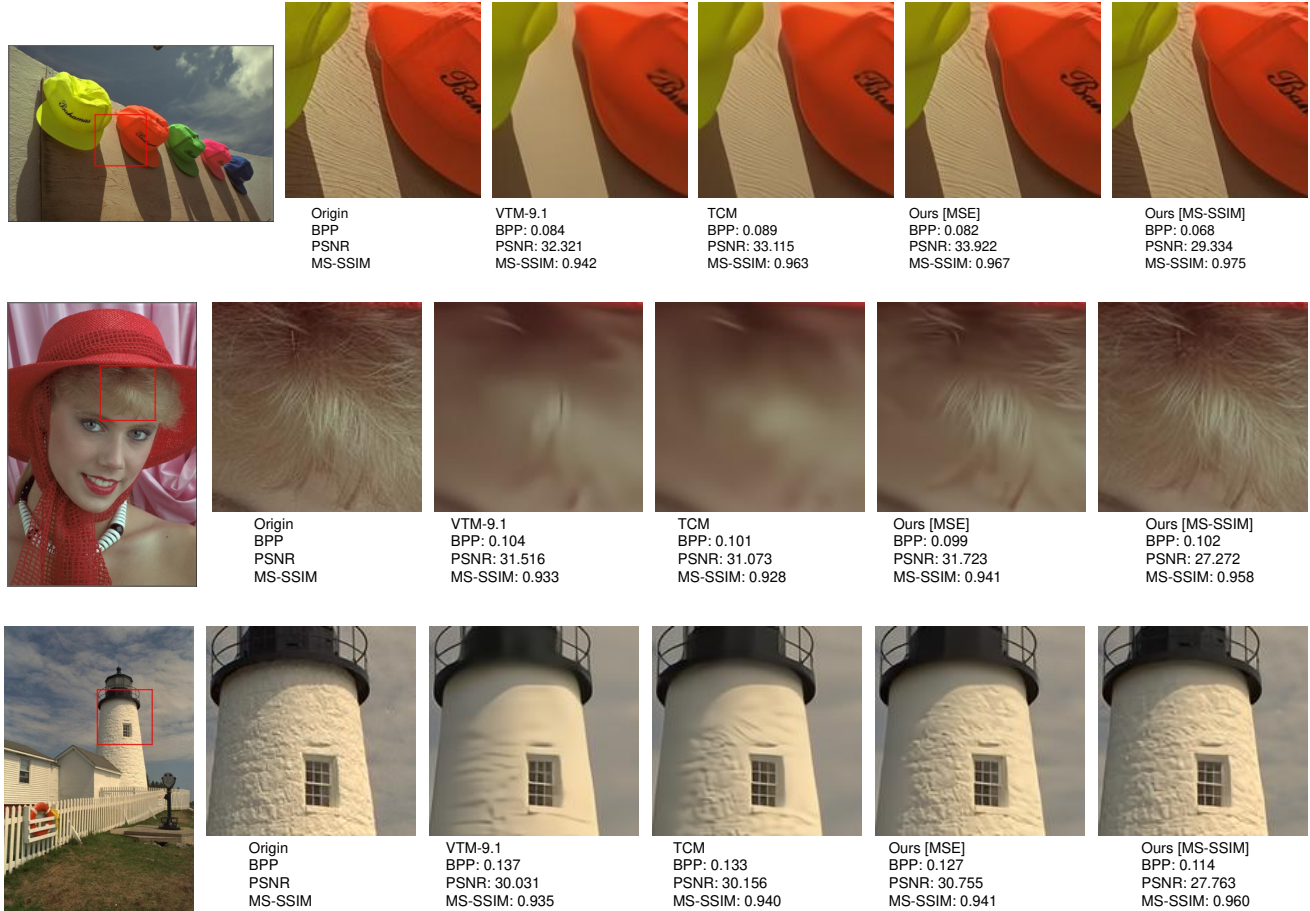


Figure 7. More Visual Comparisons.

171 patches reveals that GLIC better preserves high-frequency
 172 details and textures, effectively mitigating the severe detail
 173 loss and over-smoothing observed in other methods. This
 174 improvement stems from our complexity-aware neighbor
 175 quota assignment, which adaptively allocates more connec-
 176 tions in complex regions, leading to better compression and
 177 reconstruction of intricate structures.

178 12. Why RMS-Gradient

179 As shown in Fig. 8, for LIC models, both the latent fea-
 180 tures and the Sobel-gradient responses exhibit strong en-
 181 ergy compaction: a small subset of channels concentrates
 182 most of the energy (carrying the most critical information),
 183 while the majority of channels remain low-energy. Since
 184 RMS emphasizes larger magnitudes, it is better suited than
 185 the plain averaging. RMS highlights the high-energy, in-
 186 formative channels instead of being diluted by numerous
 187 low-energy ones. This is also empirically validated by our
 188 experimental results in main manuscript Tab. 4.

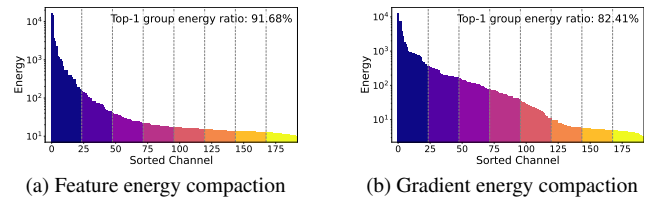


Figure 8. Channel-wise energy distribution of Stage-2 features, averaged over Kodak images. The channels are sorted into eight groups.

189 13. Limitations and future work

190 A key limitation is that we have not evaluated our dynamic
 191 GNNs' effectiveness in exploiting cross-frame redundancy
 192 for video compression. In future work, we will extend our
 193 GNN to video codecs and investigate how adaptive graphs
 194 can better model temporal correlations.

195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251**References**

- [1] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja. Compressai: a pytorch library and evaluation platform for end-to-end compression research. *arXiv e-prints*, pages arXiv–2011, 2020. 2
- [2] Jianle Chen, Yan Ye, and S Kim. Algorithm description for versatile video coding and test model 1 (vtm 1). *Joint Video Experts Team (JVET) of ITU-T SG*, 16:3–12, 2020. 4
- [3] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11030–11039, 2020. 1
- [4] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017. 1
- [5] Donghui Feng, Zhengxue Cheng, Shen Wang, Ronghua Wu, Hongwei Hu, Guo Lu, and Li Song. Linear attention modeling for learned image compression. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7623–7632, 2025. 2
- [6] Haisheng Fu, Feng Liang, Jie Liang, Yongqiang Wang, Zhenman Fang, Guohe Zhang, and Jingning Han. Fast and high-performance learned image compression with improved checkerboard context model, deformable residual module, and knowledge distillation. *IEEE Transactions on Image Processing*, 2024. 1
- [7] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6185–6194, 2023. 1
- [8] 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*, pages 5718–5727, 2022. 2, 4
- [9] Daowen Li, Yingming Li, Heming Sun, and Lu Yu. Deep image compression based on multi-scale deformable convolution. *Journal of Visual Communication and Image Representation*, 87:103573, 2022. 1
- [10] Han Li, Shaohui Li, Wenrui Dai, Chenglin Li, Junni Zou, and Hongkai Xiong. Frequency-aware transformer for learned image compression. In *The Twelfth International Conference on Learning Representations*, 2024. 4
- [11] Han Li, Shaohui Li, Wenrui Dai, Maida Cao, Nuowen Kan, Chenglin Li, Junni Zou, and Hongkai Xiong. On disentangled training for nonlinear transform in learned image compression. In *The Thirteenth International Conference on Learning Representations*, 2025. 4
- [12] Jinming Liu, Heming Sun, and Jiro Katto. Learned image compression with mixed transformer-cnn architectures. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14388–14397, 2023. 4
- [13] 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*, pages 10012–10022, 2021. 1
- [14] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *Advances in neural information processing systems*, 29, 2016. 1
- [15] Gabriele Spadaro, Alberto Presta, Enzo Tartaglione, Jhony H Giraldo, Marco Grangetto, and Attilio Fiandrrotti. Gabic: Graph-based attention block for image compression. In *2024 IEEE International Conference on Image Processing (ICIP)*, pages 1802–1808. IEEE, 2024. 1, 2
- [16] Zhisen Tang, Hanli Wang, Xiaokai Yi, Yun Zhang, Sam Kwong, and C-C Jay Kuo. Joint graph attention and asymmetric convolutional neural network for deep image compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(1):421–433, 2022. 1, 2
- [17] Chunhui Yang, Yi Ma, Jiayu Yang, Shiyi Liu, and Ronggang Wang. Graph-convolution network for image compression. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 2094–2098. IEEE, 2021. 1, 2
- [18] Fanhu Zeng, Hao Tang, Yihua Shao, Siyu Chen, Ling Shao, and Yan Wang. Mambaic: State space models for high-performance learned image compression. *arXiv preprint arXiv:2503.12461*, 2025. 4