# WINDOW-BASED HIERARCHICAL DYNAMIC ATTENTION FOR LEARNED IMAGE COMPRESSION

Anonymous authors

Paper under double-blind review

# Abstract

Transformers have been successfully applied to learned image compression (LIC). In fact, dense self-attention is difficult to ignore contextual information that degrades the entropy estimations. To overcome this challenging problem, we incorporate dynamic attention in LIC for the first time. The window-based dynamic attention (WDA) module is proposed to adaptively tune attention based on the entropy distribution by sparsifying the attention matrix. Additionally, the WDA module is embedded into encoder and decoder transformation layers to refine attention in multi-scales, hierarchically extracting compact latent representations. Similarly, we propose the dynamic-reference entropy model (DREM) to adaptively select context information. This decreases the difficulty of entropy estimation by leveraging the relevant subset of decoded symbols, achieving an accurate entropy model. To the best of our knowledge, this is the first work employing dynamic attention for LIC. Extensive experiments demonstrate the proposed method outperforms the state-of-the-art LIC methods.

027

004

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

028 Vision Transformer (ViT) has achieved tremendous advancements in the field of computer vision, 029 with many studies applying it to learned image compression (LIC) methods (Zhu et al., 2022; Qian et al., 2022a;b; Liu et al., 2023). Efficient self-attention between all sequence elements helps the model pay attention to long-range information. There are mainly two aspects of works transferring 031 CNN-based learned image compressions to ViT architectures. Utilizing Swin Transformers (Swin-T) (Liu et al., 2021b) in main encoder-decoders to build powerful nonlinear transforms (Liu et al., 033 2023; Zhu et al., 2022; Zou et al., 2022; Lu et al., 2022). On the other hand, some works leverage 034 ViTs in entropy models to capture global contextual information, supporting a more accurate probability estimation of the latent representation distributions (Qian et al., 2022a; Koyuncu et al., 2022; Kim et al., 2022). However, the rate-distortion (RD) performance improvements of these methods 037 are marginal.

038 Nonlinear transformations and entropy models are the key components of LIC. Despite ViTs enables attention to more distant context, it does not guarantee compact transformations and accurate entropy 040 estimation. Adjacent features exhibit stronger causal relationships and the previous work (Minnen 041 et al., 2018) reveals that convolution kernel sizes larger than  $5 \times 5$  (with larger receptive fields) 042 unexpectedly compromise the RD performance. The works (He et al., 2021; Zou et al., 2022) also 043 prove that redundancy primarily exists in local regions. Some redundancy information indeed exists 044 in distant regions (Qian et al., 2022b), but referring global contextual information increases the risk of overfitting. Previous works have focused solely on the long-range modeling capabilities of ViTs, ignoring the issue of overfitting. 046

In this paper, we analyze the challenge of applying ViTs to image compression, and a novel method is proposed: dynamically sparsifying attention. Plain Swin-T compute paired attention between all elements in local windows. In other works, all decoded symbols in a window are considered when decoding the current node. Different from recognition tasks, the goal of compression is to remove redundancy. Reference to irrelevant content can mislead probability estimations. We claim that the core contradiction of overfitting is that the attention-pattern space of ViTs is much large than redundancy-pattern space. To overcome this problem, we sparsify the attention-pattern space and propose a compression model adopts adaptive attention patterns learning from the entropy distribu-

tion of the image. Specifically, the model is built on the window-based dynamic attention (WDA)
 module, which tunes the attention matrix to ignore useless references in local windows. The WDA
 module works in multiple feature scales to hierarchically tune long-range attention patterns. Further more, the dynamic-reference entropy model (DREM) is proposed, which builds upon the concept
 of dynamic attention patterns, adaptively selecting reference contextual information for the current
 encoding element based on the known entropy distribution. The aggregation of relevant decoded
 symbol subsets significantly reduces the difficulty of probability estimation in the entropy model.

To the best of our knowledge, we first focus on the overfitting issue of transformer-based learned
 image compression methods and modulate the attention by the means of dynamic sparsification
 patterns. In summary, our contributions can be concluded as follows:

- We first integrate dynamic attention into learned image compression, which narrows the gap of attention space and redundancy space. Sparse attention patterns ignore globally irrelevant contexts, reducing the risk of overfitting.
- We propose the window-based dynamic attention (WDA) module, which adaptively learns attention patterns from entropy information of latent representations. The WDA modulates the attention matrix at different scales in a fine-to-coarse manner.
- We present the dynamic-reference entropy model (DREM) to select subsets of decoded symbols, which provide enough contextual information and simultaneously reduce the optimization difficulty.
  - Experiment results demonstrate that our proposed method achieves 13.42%, 17.74% and 12.93% BD-rate gains over VTM-17.0 on the Kodak, Tecnick and CLIC datasets respectively and outperforms the state-of-the-art LIC method MLIC++.
- 2 RELATED WORK
- 079 080 081

065

066

067 068

069

070 071

073

074

075

076

077

2.1 LEARNED IMAGE COMPRESSION

Early LIC methods adopt convolutional neural networks (CNNs) in both encoder-decoders and en-083 tropy models (Ballé et al., 2018; Minnen et al., 2018; Lee et al., 2018). (Cheng et al., 2020) first 084 incorporates the attention mechanism into LIC, which pays more attention to regions with com-085 plicated textures. However, the local receptive field of CNNs limits their ability to capture longrange spatial dependencies. Some global methods utilize non-local networks (Chen et al., 2021) and 087 content-weighted attention masks (Li et al., 2018; Mentzer et al., 2018) to alloacte bits across the 088 entire image, leading to an overall improvement in RD performance. With the rise of transformers, 089 ViTs are gradually emerging in LIC. The global self-attention constructs more powerful nonlinear transformations (Lu et al., 2022; Zhu et al., 2022; Liu et al., 2023; Zou et al., 2022) and provide rich contextual information in entropy models (Qian et al., 2022a; Kim et al., 2022; Liu et al., 2023). 091 However, the downside of global information is that irrelevant context increases the difficulty of 092 entropy estimation and the risk of overfitting. We propose the WDA module to dynamically sparsify 093 the attention patterns to address the problem. 094

095

096 2.2 DYNAMIC ATTENTION

Previous works have demonstrated that a significant amount of computational redundancy exists 098 in ViTs. Only a small proportion of tokens contribute to the final prediction, thus removing those useless tokens improves the computational efficiency without harming the performance (Chen et al., 100 2023; Wei et al., 2023). Following that, some sparse attention methods are proposed to accelerate 101 ViTs, including token sampling (Rao et al., 2021; Fayyaz et al., 2022; Tang et al., 2022) and atten-102 tion masking (Liu et al., 2021a; Kitaev et al., 2019). Among those methods, static sparse methods 103 (Tay et al., 2020; Kong et al., 2022) introduce heuristic sparse attention patterns with challenging of 104 generalization. While dynamic methods (Yin et al., 2022; Venkataramanan et al., 2023; Lee et al., 105 2024) learn dynamic attention patterns from data in a flexible way. Our work is inspired by dynamic sparse attention but applied in a different domain. Specifically, we apply the dynamic sparse at-106 tention to more efficiently eliminate representation redundancy rather than to reduce computational 107 complexity.



Figure 1: Overall framework of the proposed Window-based Dynamic Attention Learned Image Compression (WDA-LIC).  $g_a$  and  $g_s$  denote analysis and synthesis transforms consist of multiple Conv Blocks and WDA modules. Conv Blocks extract local features and the WDA modules capture long-range contextual information with dynamic attention patterns. N denotes output channel numbers in every layers.

134

# 2.3 CONTEXT ENTROPY MODELING

135 In LIC entropy models replace the marginal probability distribution of latent variables by the joint 136 probability distribution with prior variables to reduce entropy (Ballé et al., 2018; Minnen et al., 2018; Lee et al., 2018). Due to the sequential property of decoding, leveraging previously decoded fea-137 tures (*i.e.*, context) to provide predictive information for the current decoding step can significantly 138 reduce the joint entropy. And an optimal contextual pattern determines the upper bound of predic-139 tion accuracy. Some works divide feature channels into multiple slices and remove redundacy by 140 leveraging correlations between channels (Minnen & Singh, 2020; He et al., 2022; Zhu et al., 2022). 141 In terms of spatial redundancy, CNN-based methods capture local correlations between neighboring 142 representations (He et al., 2021; Zou et al., 2022; Guo et al., 2021) and transformer-based methods 143 calculate relevant information over longer ranges (Liu et al., 2023; Lu et al., 2022). Although rele-144 vant information may exist in global regions, previous works (He et al., 2021; Minnen et al., 2018) 145 show that adjacent pixels are likely to have a stronger causal relationship. Focusing on too much 146 irrelevant information increases the difficulty of prediction. Some methods select top-K elements in global regions to centralize attention (Qian et al., 2022a;b; Ma et al., 2021). However, this fixed-147 number reference pattern fails to adapt to sample differences. We propose the dynamic-reference 148 entropy model (DREM) to adaptively select reference subsets with entropy information. 149

### 150 151

152

154

158

# 3 Methods

# 153 3.1 PROBLEM FORMULATION

The architecture of our proposed Window-based Hierarchical Dynamic Attention Learned Image
 Compression (WDA-LIC) is shown in Figure 11. The overall algorithmic can be formulated as
 follows:

$$y = g_a(x; \theta_{g_a}), \hat{y} = Q(y), \hat{x} = g_s(\hat{y}; \theta_{g_s}),$$
(1)

159 where the encoder  $g_a$  with parameters  $\theta_{g_a}$  transforms the input image x to latent representation y. 160 Following that y is quantized to  $\hat{y}$ , which is modeled as a single Gaussian distribution with esti-161 mated parameters  $(\mu, \sigma)$  to be entropy encoded. The decoder  $g_s$  with parameters  $\theta_{g_s}$  utilizes  $\hat{y}$  to reconstruct  $\hat{x}$ . It is so critical to accurately estimate the distribution parameters  $\mu$  and  $\sigma$ . We adopt 162 the hyperprior model (Ballé et al., 2018) and the channel-wise autoregression entropy model (Min-163 nen & Singh, 2020; He et al., 2022) to estimate the Gaussian parameters ( $\mu, \sigma$ ). The hyperprior 164 model is used to capture side information: 165

174 175

176 177

178 179

183

185 186 187

189

190 191

192

199

200

 $z = h_a(y; \phi_{h_a}), \hat{z} = Q(z), \psi_h = h_s(\hat{z}; \phi_{q_s}),$ (2)

167 where  $\psi_h$  denotes the side information provided by the hyperprior model. A mount of redundancy 168 exists between channels and the decoding process is sequential, we follow provious works (He et al., 169 2022; Jiang et al., 2023) to divided latent variables y into S slices  $\{y^0, y^1, \dots, y^{s-1}\}$  so that encoded 170 slices provide contextual information to help the entropy estimation of currently encoding slice as shown in Figure 3. During the process of encoding slice  $y_i$ , all its front slices  $\{\hat{y}^0, \hat{y}^1, \dots, \hat{y}^{i-1}\}$ 171 and the side information  $\psi_h$  are fed into the proposed Dynamic-Reference Entropy Model (DREM) 172 to estimate the Gaussian parameter of current slice as follows: 173

$$\Phi_i = e(\psi_h, \hat{y}^{< i}, y^i) 
= (\mu_i, \sigma_i),$$
(3)

where *e* is the DREM. Therefore, the probability of current slice is considered as follows:

$$p_{\hat{y}^{i}|\hat{z},\hat{y}^{< i}}(\hat{y}^{i} \mid \hat{z},\hat{y}^{< i}) \sim \mathcal{N}(\mu^{i},\sigma^{i}),$$
(4)

181 Since the entropy bottleneck  $\Psi$  is used to encode  $\hat{z}$  as  $p_{\hat{z}|\Psi}(\hat{z} \mid \Psi)$ , the overall rate-distortion (RD) 182 loss function is defined as:

$$\mathcal{L} = \mathcal{R}(\hat{y}) + \mathcal{R}(\hat{z}) + \lambda \cdot \mathcal{D}(x, \hat{x})$$
  
=  $\mathbf{E}[-\log_2(p_{\hat{y}\mid\hat{z}}(\hat{y}\mid\hat{z}))] + \mathbf{E}[-\log_2(p_{\hat{z}\mid\Psi}(\hat{z}\mid\Psi))]$   
+  $\lambda \cdot \mathcal{D}(x, \hat{x}),$  (5)

where  $\lambda$  is a Lagrangian multiplier to control the RD tradeoff.  $\mathcal{D}(x, \hat{x})$  denotes the distortion term 188 such as Mean squared error (MSE) loss.  $\mathcal{R}(\hat{y})$  and  $\mathcal{R}(\hat{z})$  are the bit rates of latent representations  $\hat{y}$ and  $\hat{z}$ .

### 3.2 DYNAMIC ATTENTION-BASED TRANSFORMATION

193 We build the nonlinear transformations in a CNN-Transformer mixed way, as shown in Figure . At 194 each stage, the Conv Block extracts features through a CNN, followed by a Swin-T based WDA module to fuse features within a local window. Specifically, the WDA module adopts adaptive attention patterns with the instruction of entropy distribution in local windows. With smaller feature 196 scales, the WDA module tunes the attention locations in a larger receptive field during the trans-197 forming. The following sections elaborate our proposed WDA module.

#### WINDOW-BASED DYNAMIC ATTENTION MODULE 3.2.1

201 The WDA module is illustrated in Figure 2 and can be viewed as a Swin-T with dynamic attention patterns. Given an input feature  $X \in \mathbb{R}^{H \times W \times C}$ , the vanilla Swin-T divides X into non-overlapping 202 partitions  $[X_1, \ldots, X_M]$  with  $K \times K$  size windows or shifted-windows and arrange them into the 203 feature matrix, where  $\mathbf{X}_i \in \mathbb{R}^{N \times C}$ ,  $N = K \times K$ ,  $1 \le i \le M$  and  $M = \frac{H}{K} \times \frac{W}{K}$ . The multi-head 204 205 self-attention is conducted within each window  $X_i$  as follows: 206

207

$$egin{aligned} oldsymbol{Q}, oldsymbol{K}, oldsymbol{V} &= oldsymbol{X}_i oldsymbol{W}^{oldsymbol{Q}}, oldsymbol{X}_i oldsymbol{W}^{oldsymbol{K}}, oldsymbol{Z}_i oldsymbol{W}^{oldsymbol{K}}, oldsymbo$$

(6)

212 213

where  $\mathbf{W}^{\mathbf{Q}}, \mathbf{W}^{\mathbf{K}}, \mathbf{W}^{\mathbf{V}} \in \mathbb{R}^{C \times d_k}$  are learnable parameters and  $d_k$  is the intermediate feature di-214 mension. The above formula calculates the self-attention of each token with all other tokens in the 215 sequence and the final output is a weighted average of all tokens within the window. Obviously,



Figure 2: Proposed window-based dynamic attention module. The attention pattern is a 0-1 matrix with the same size of attention matrix.  $\otimes$  denotes the matrix multiplication operation and  $\odot$  denotes the Hadamard dot multiplication operation.

237 spatial structure exists in attention and awful context spreads the distribution of latent features, in-238 creasing the entropy. This is not supposed to happen when coding. The correct approach is to focus 239 attention on tokens with rich mutual information, which makes the latent features more compact and reduces entropy. Though some previous methods (Qian et al., 2022a;b) utilize Top-K scheme to 240 select K-most relevant reference elements, this fixed attention pattern cannot adapt to the diversity 241 of image distributions, thus the improvement of RD performance is marginal. Intuitively, regions 242 with complex texture should reference more contextual information. Our WDA module adopt dy-243 namic attention patterns. For each window, we can easily obtain the covariance matrix between 244 latent variables before computing the attention matrix V as follows: 245

246 247

236

248 249

250

251

257

258

259 260 261

262

263

264 265 266  $V = \frac{1}{c-1} (X_i - \mu) (X_i - \mu)^T,$ with  $\mu = \frac{1}{C} \sum_{j=1}^C X_i^j,$  (7)

where  $V \in \mathbb{R}^{N \times N}$ ,  $\mu$  is the mean of each token, *C* is the number of channels. The diagonal elements of *V* represent the variance of each variable, where larger variances correspond to higher entropy. The other elements indicate the correlation between pairs of variables. Although the covariance matrix only represents linear correlations, it is sufficient as a clue for attention aggregation. To dynamically sparsify the attention matrix, we obtain the mask matrix *M* as follows:

 $\boldsymbol{M}(i,j) = \begin{cases} 0 & \text{if } \mid \frac{\boldsymbol{V}(i,j)}{\boldsymbol{V}(i,i)} \mid \geq t, \\ -inf & otherwise, \end{cases}$ (8)

where t is the threshold, and is set to be 0.8. It is obvious that the attention pattern is dynamic due to the diversity of entropy (*i.e.*, V(i, j)). Following that the mask matrix modulates the attention matrix as follows:

$$\hat{A}(X_i) = \operatorname{softmax}(\frac{Q \cdot K^T}{\sqrt{d_k}} + M),$$
(9)

It is equal to adopt the Hadamard operation with 0 - 1 masks shown in Figure 2. And the final output of the WDA module is as:

$$\hat{O}(X_i) = \hat{A} \cdot V, \tag{10}$$

#### 270 3.2.2 HIERARCHICAL ATTENTION MODULATION 271

272 As the features are downsampled, the receptive field of each window corresponding to the original 273 image gradually increases. To modulate the attention pattern in a hierarchical way, we apply the WDA module to multiple feature scales. When tuning attention at features with d-downsampled 274 scales, the receptive filed to the image can be expressed as: 275

$$\mathcal{F} = \left(\frac{3}{2}K \times \frac{3}{2}K \times d\right)^2,\tag{11}$$

where  $\mathcal{F}$  denotes the resolution of the original image K is the window size. The factor  $\frac{3}{2}$  is due to the shifting-window operation. Larger K have a larger receptive field with more irrelevant tokens, thus the threshold t tends to increase to abandon those tokens.

### 3.3 DYNAMIC-REFERENCE ENTROPY MODEL

284 Figure shows the pipeline of the Dynamic-Reference En-285 tropy Model (DREM). To encode the latent slice  $y^i$ , all 286 previously encoded slices  $\hat{y}^{< i}$  and hyperprior context  $\psi_h$ 287 are utilized. Specifically, we split  $y^i$  into two parts (*i.e.*, 288  $y_a^i$  and  $y_n^i a$  in the checkerboard spatial pattern following 289 the previous works (He et al., 2021; Jiang et al., 2023). After coding  $\hat{y}_a^i$ , it provides local spatial information to 290  $y_n^i a$ . Channel-wise and global spatial context are pre-291 dicted from previously encoded slices  $\hat{y}^{< i}$ . All context 292 representations are concatenated in channel dimension 293 and fed into the entropy estimation network to predict the distribution parameters  $(\mu, \sigma)$ . 295

The dynamic-reference is reflected in the selection of 296 global contextual information. As equation 8, we lever-297 age the WDA module in the global context network to 298 dynamically select a subset of tokens in  $\hat{y}^{<i}$  according to 299 the current entropy distribution of  $y^i$ . The workflow of 300 DREM can be summarized as follows: 301

$$\begin{split} \Phi_{ch}^{i} &= g_{ch}(\hat{y}^{(12)$$

306 307 308

309 310

311 312

313

314

315

316

317

318

319

302

303

305

276 277 278

279

280

281 282

283

#### 4 EXPERIMENTS

## 4.1 EXPERIMENTAL SETUP



320 321 322

323

## 4.2 RATE-DISTORTION PERFORMANCE



Figure 3: Pipeline of the Dynamic-Reference Entropy Model (DREM). where  $g_{ch}, g_{ql}, g_{lc}, g_{en}$  are networks and  $\Phi_{ch}^{i}, \Phi_{al}^{i}, \Phi_{lc}^{i}, \Phi_{a}^{i}, \Phi_{na}^{i}$  denote context latent variables. More network architecture details can be found in Appendix A.

<sup>&</sup>lt;sup>1</sup>https://interdigitalinc.github.io/CompressAI/



Figure 4: RD performance on the Kodak dataset (left:PSNR, right:MS-SSIM).

343 Six popular LIC methods (Jiang et al., 2023; Liu 344 et al., 2023; Zou et al., 2022; He et al., 2022; Cheng 345 et al., 2020; Ballé et al., 2018) and the traditional 346 codec VTM-17.0 is compared. The R-D perfor-347 mance of the Kodak dataset is shown in Figure 4. 348 The results of CLIC and Tecnick datasets are presented in Figure in the Appendix. To comprehen-349 sively compare the RD performance of two comres-350 sion methods, we utilize the BD-Rate (Bjontegaard, 351 2001) metric. 352

353 354

355

356

340 341 342

### 4.3 ABLATION STUDY

357 Effectiveness of the WDA module. We remove the 358 WDA module in analysis and synthesis transformations as the baseline and compare the BD-rate over 359 VTM-17.0. The results are displayed in Table 2, 360 which illustrates the efficiency of the WDA module. 361 Atten denotes plain attention patterns that discards 362 masks. w/ WDAtten (n=1) represents the method 363 that maintains the last WDA module and w/ WDAt-364 ten (n=4) maintains all WDA modules. The WDA module is lightweight and is easy to be compatible 366 with other networks. The results further shows that 367 retaining the last WDA module still keeps perfor-368 mance advantage. The visualization of latent distri-

Table 1: BD-rate results over VTM-17.0 of state-of-the-art LICs. The evaluation is conducted on the Kodak dataset.

Methods	PSNR	MS-SSIM
VTM-17.0	-	-
Cheng(CVPR2020)	+5.58	-44.21
He(CVPR2022)	-5.59	-44.60
Zou(CVPR2020)	-2.48	-47.72
Liu(CVPR2023)	-10.14	-48.94
Jiang(ACMMM2023)	-13.39	-53.63
Ours	-13.42	-53.96

Table 2: BD-rate results over VTM-17.0 on the CLIC Valid dataset of different models.

Methods	$\operatorname{Params}(M)$	BD-rate
w/o Atten	50.34	-9.08
w/ Atten	60.48	-11.64
w/ WDAtten (n=1)	52.75	-12.08
w/ WDAtten (n=4)	60.48	-12.93
VTM-17.0	-	0

butions are shown in Figure 5. It is obvious that the WDA module compacts the distribution of latent representations.

Performance of DREM. DREM is proposed to dynamically selecting reference subsets of tokens
in global range. To illustrate the performance of DREM, we compare our proposed method with
different global attention patterns. The global context network is abandoned to build the baseline.
The highlight of DREM is adaptivity. Therefore we compare with fixed attention patterns as shown
in Table 3. To illustrate the importance of global context informarion, the global context network
is discarded. The fully method computes pairwise correlations without attention masks. The Top-K
method maintain K most relevant tokens in each prediction and the number K is chosen empirically (Qian et al., 2022a;b), lacking of flexibility.

381		Methods	BPP	PSNR			
382			0.211	20.975			
383		$W/0 g_{gl}$	0.311	30.875			
384		Fully	0.279	32.118			
385		Тор-К	0.288	31.980			
386		DREM	0.271	32.376			
387							
388							
389		1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1		CONTRACTOR STORE	20	1 C 1 1 1 1	
390			20	10 C		and the second	10
391			-0	ALC: NO	0	Contraction of the second	0 -10
392			-20	CALLER CO.	-20	1110	-20
393			200	AAAA		N 1 1 1	-30
394	Kodim03	v(w/o Atten)		v(w/ Atten)		v(w/WDAtten)	
395	Rounios	y(w/o / titen)	Le Contraction	y(w/ Attell)	6	y(w, w brace)	6
396			4	A DECK OF A DECK	-4		•4
397		ALC: NOT THE REAL PROPERTY OF	·2		2		•2
398		State State	-2		·0	Star Links	•0
399		States and a state of	-4	States and states	Ľ	ALC: NO. OF THE OWNER.	-2
400		$\sigma$ (w/o Atten)		$\sigma$ (w/ Atten)		$\sigma$ (w/ WDAtten)	
401			40			1000	20
402			20	18 C	20		10
403			0	22.007	•0		0
404			-20		-20		-20
405			-40			1000	-30
406	Kodim15	y(w/o Atten)		y(w/ Atten)		y(w/ WDAtten)	
407		1. 化化学学 化化学		a lotter an art	4	A CONTRACTOR	4
408		- 136 9 GMG	2.5	A STATE OF STREET	2	States and	•2
409			0.0	State State State	0	Sim State of the	10
410		and the second second	2.5	State of the local division of the local div	-2		·-2
411		STATE STATE STATE	-5.0	State of the second	-4	COMPAREMENTS	-4
412		$\sigma$ (w/o Atten)		$\sigma(\text{w/Atten})$		$\sigma$ (w/ WDAtten)	
413							

378 Table 3: Comparison of different attention patterns. The RD perfermance on Kodak datasets are 379 displayed.

Figure 5: The average scaled deviation  $\sigma$  and feature y across channels. The model (w/o Atten) abandons the WDA module and (w/ Atten) and (w/ WDAtten) denote adopt vanilla attention patterns and dynamic sparse attention patterns with the WDA module respectively.

414

415

416

380

# 419

5

CONCLUSION

- 420
- 421

# 422

424

In this paper we first adopt dynamic attention into learned image compression. Based on the as-423 sumption that the redundancy information densely distributes in local regions and sparsely exists in long-range distance, we propose the WDA module to dynamically sparsifying the attention matrix 425 in Swin-T blocks, making adaptive attention patterns learned from data possible. This is reasonable 426 because of the diversity of image entropy distribution. The WDA module dynamically modulate 427 the contextual information according to the local entropy, where regions with large entropy could 428 be allocated more long-range context. Appling the WDA into the entropy model, the proposed dynamic-reference entropy model select a subset of reference tokens, sparsifing the optimization 429 space and decreases the risk of overfitting. Extensive experiments demonstrate the performance ad-430 vantage of our method and proves the possibility for compression networks to evolve in a dynamic 431 and flexible direction in the future.

# 432 REFERENCES

- Nicola Asuni, Andrea Giachetti, et al. Testimages: a large-scale archive for testing visual devices and basic image processing algorithms. In *STAG*, pp. 63–70, 2014.
- 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.
- 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.
- 443 Gisle Bjontegaard. Calculation of average psnr differences between rd-curves. In VCEG-M33, 2001.
- 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*, 30:3179–3191, 2021.
- Xuanyao Chen, Zhijian Liu, Haotian Tang, Li Yi, Hang Zhao, and Song Han. Sparsevit: Revisiting activation sparsity for efficient high-resolution vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2061–2070, 2023.
- 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*, pp. 7939–7948, 2020.
- Mohsen Fayyaz, Soroush Abbasi Koohpayegani, Farnoush Rezaei Jafari, Sunando Sengupta, Hamid
  Reza Vaezi Joze, Eric Sommerlade, Hamed Pirsiavash, and Jürgen Gall. Adaptive token sampling
  for efficient vision transformers. In *European Conference on Computer Vision*, pp. 396–414.
  Springer, 2022.
- Zongyu Guo, Zhizheng Zhang, Runsen Feng, and Zhibo Chen. Causal contextual prediction for
   learned image compression. *IEEE Transactions on Circuits and Systems for Video Technology*,
   32(4):2329–2341, 2021.
- 462 Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. Checkerboard context
  463 model for efficient learned image compression. In *Proceedings of the IEEE/CVF Conference on*464 *Computer Vision and Pattern Recognition*, pp. 14771–14780, 2021.
- 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*, pp. 5718–5727, 2022.
- Wei Jiang, Jiayu Yang, Yongqi Zhai, Peirong Ning, Feng Gao, and Ronggang Wang. Mlic: Multireference entropy model for learned image compression. In *Proceedings of the 31st ACM Inter- national Conference on Multimedia*, pp. 7618–7627, 2023.
- 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*, pp. 5992–6001, 2022.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2019.
- 479 Eastman Kodak. Kodak lossless true color image suite (photocd pcd0992). 1993.
- Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *European conference on computer vision*, pp. 620–640. Springer, 2022.
- 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 *European Conference on Computer Vision*, pp. 447–463. Springer, 2022.

504

- Heejun Lee, Jina Kim, Jeffrey Willette, and Sung Ju Hwang. Sea: Sparse linear attention with
   estimated attention mask. In *The Twelfth International Conference on Learning Representations*, 2024.
- Jooyoung Lee, Seunghyun Cho, and Seung-Kwon Beack. Context-adaptive entropy model for endto-end optimized image compression. In *International Conference on Learning Representations*, 2018.
- Mu Li, Wangmeng Zuo, Shuhang Gu, Debin Zhao, and David Zhang. Learning convolutional
   networks for content-weighted image compression. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3214–3223, 2018.
- Jiaheng Liu, Guo Lu, Zhihao Hu, and Dong Xu. A unified end-to-end framework for efficient deep image compression. *arXiv e-prints*, pp. arXiv–2002, 2020.
- Jinming Liu, Heming Sun, and Jiro Katto. Learned image compression with mixed transformer cnn architectures. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition
   (CVPR), pp. 14388–14397. IEEE Computer Society, 2023.
- Liu Liu, Zheng Qu, Zhaodong Chen, Yufei Ding, and Yuan Xie. Transformer acceleration with
   dynamic sparse attention, 2021a. URL https://arxiv.org/abs/2110.11299.
- 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 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9992–10002. IEEE, 2021b.
- Ming Lu, Peiyao Guo, Huiqing Shi, Chuntong Cao, and Zhan Ma. Transformer-based image compression. In 2022 Data Compression Conference (DCC), pp. 469–469. IEEE, 2022.
- Changyue Ma, Zhao Wang, Ruling Liao, and Yan Ye. A cross channel context model for latents in deep image compression. *arXiv e-prints*, pp. arXiv–2103, 2021.
- Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Condi tional probability models for deep image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4394–4402, 2018.
- 516 David Minnen and Saurabh Singh. Channel-wise autoregressive entropy models for learned image
   517 compression. In 2020 IEEE International Conference on Image Processing (ICIP), pp. 3339– 3343. IEEE, 2020.
- David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for learned image compression. *Advances in neural information processing systems*, 31, 2018.
- Yichen Qian, Xiuyu Sun, Ming Lin, Zhiyu Tan, and Rong Jin. Entroformer: A transformer-based
   entropy model for learned image compression. In *International Conference on Learning Representations*, 2022a.
- Yichen Qian, Zhiyu Tan, Xiuyu Sun, Ming Lin, Dongyang Li, Zhenhong Sun, Li Hao, and Rong Jin.
   Learning accurate entropy model with global reference for image compression. In *International Conference on Learning Representations*, 2022b.
- Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit:
   Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021.
- Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch
   slimming for efficient vision transformers. In *Proceedings of the IEEE/CVF Conference on Com- puter Vision and Pattern Recognition*, pp. 12165–12174, 2022.
- Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In International Conference on Machine Learning, pp. 9438–9447. PMLR, 2020.
- George Toderici, Wenzhe Shi, Radu Timofte, Lucas Theis, Johannes Ballé, Eirikur Agustsson, Nick
   Johnston, and Fabian Mentzer. Clic: Workshop and challenge on learned image compression. In
   *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit*, 2020.

- Shashanka Venkataramanan, Amir Ghodrati, Yuki M Asano, Fatih Porikli, and Amir Habibian. Skipattention: Improving vision transformers by paying less attention. In *The Twelfth International Conference on Learning Representations*, 2023.
- Cong Wei, Brendan Duke, Ruowei Jiang, Parham Aarabi, Graham W Taylor, and Florian Shkurti.
   Sparsifiner: Learning sparse instance-dependent attention for efficient vision transformers. In
   *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22680–22689, 2023.
- Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-vit:
   Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10809–10818, 2022.
  - Yinhao Zhu, Yang Yang, and Taco Cohen. Transformer-based transform coding. In *International Conference on Learning Representations*, 2022.
  - Renjie Zou, Chunfeng Song, and Zhaoxiang Zhang. The devil is in the details: Window-based attention for image compression. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 17471–17480. IEEE Computer Society, 2022.

# A APPENDIX

551

552

553

554

555

556

558

559

561

### A.1 DETAILS OF THE ARCHITECTURES











