
LiteVAR: Compressing Visual Autoregressive Modelling with Efficient Attention and Quantization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Visual Autoregressive (VAR) has emerged as a promising approach in image genera-
2 tion, offering competitive potential and performance comparable to diffusion-based
3 models. However, current AR-based visual generation models require substan-
4 tial computational resources, limiting their applicability on resource-constrained
5 devices. To address this issue, we conducted analysis and identified significant
6 redundancy in three dimensions of the VAR model: (1) the attention map, (2)
7 the attention outputs when using classifier free guidance, and (3) the data preci-
8 sion. Correspondingly, we proposed efficient attention mechanism and low-bit
9 quantization method to enhance the efficiency of VAR models while maintaining
10 performance. With negligible performance lost (less than 0.056 FID increase), we
11 could achieve 85.2% reduction in attention computation, 50% reduction in overall
12 memory and 1.5x latency reduction. To ensure deployment feasibility, we devel-
13 oped efficient training-free compression techniques and analyze the deployment
14 feasibility and efficiency gain of each technique.

15 1 Introduction

16 Visual Autoregressive (VAR [12]) modeling has explored the autoregressive (AR) paradigm for visual
17 generation, achieving performance comparable to state-of-the-art diffusion models. By leveraging the
18 "multi-scale" nature of images, VAR introduces a scale-by-scale generation scheme, progressing from
19 coarse to fine. However, despite operating on high-level visual tokens, the VAR generation process still
20 requires iterative token generation across multiple scales, resulting in substantial computational cost.
21 This challenge hinders the broader application of VAR models on resource-constrained platforms,
22 highlighting the need for efficiency improvements. In this paper, we focus on designing training-free
23 model compression techniques to reduce the computational and memory burden of VAR models. We
24 hope our research could shed some lights on practical acceleration of VAR and even more AR-based
25 image generative models [11, 21].

26 Based on algorithmic characteristics, we explore the redundancy for VAR to design corresponding
27 optimization. As presented in Fig. 1, we conclude the redundancy in the following dimensions:

28 **Redundancy In Attention Map.** As discussed in prior literature on vision transformer [5], visual
29 models tend to exhibit a local feature extraction nature. Using global attention that aggregates all
30 tokens may therefore be redundant, with much of the computation spent on representing relatively
31 weak long-range relationships between visual tokens. Inspired by this, we visualize the attention
32 map of the VAR model in Fig. 2-(a) and find that tokens primarily focus on their local window in
33 the attention map, while most attention values for distant tokens are close to zero. Additionally, we
34 observe a unique "multi-diagonal" pattern in the VAR attention map, where visual tokens are locally
35 aggregated within each scale.

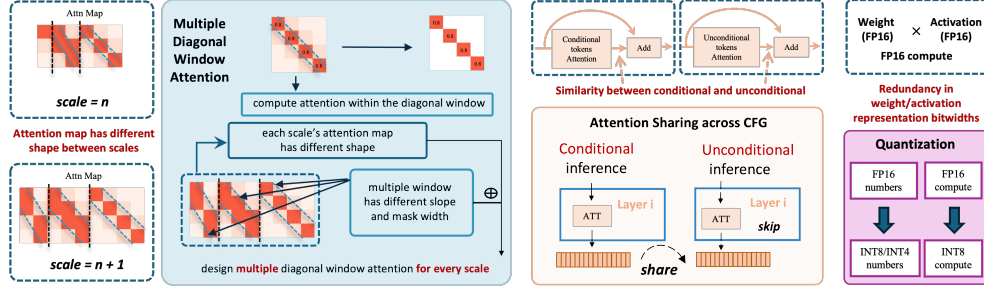


Figure 1: **Three dimensions of redundancy and corresponding compression techniques.** We discover redundancy exists in the attention map level, the classifier free guidance level, and the representation data precision level. We design the multi-diagonal windowed attention, CFG-wise sharing, and mixed precision quantization to address the above redundancy.

36 In order to leverage the unique characteristics of VAR attention maps, we propose replacing global
 37 attention with windowed local attention at each stage, which we term “multidiagonal windowed
 38 attention”. This approach effectively reduces both the computational and memory costs of attention.
 39 By incorporating multi-diagonal windowed attention, we could save 70-80% of attention computation
 40 without compromising performance. While attention computation is not a critical bottleneck in the
 41 current experimental setting (VAR on ImageNet 256x256) due to the relatively low resolution, it is
 42 important to note that attention costs scale quadratically with token length. A recent research [17]
 43 suggests that for 2K resolution generation, attention computation can become the primary bottleneck.

44 **Redundancy In Attention Outputs When Using Classifier-Free Guidance (CFG).** The CFG
 45 technique [2] is widely applied in conditioned generation, not only for diffusion models but also
 46 for autoregressive (AR) models [6, 12]. In this technique, the model is run twice—once with
 47 and once without the control signal—and the outputs are combined via a weighted sum. The
 48 weighting coefficient controls the strength of the control signal. Recent studies [17] have identified
 49 computational redundancy between the conditional and unconditional inferences in diffusion models.
 50 In this work, we investigate whether similar redundancy exists in AR models, using VAR as a
 51 representative example. By visualizing the similarity between the attention QKV of the conditional
 52 and unconditional branches in VAR generation (in Fig. 2)-(b), we observed significant overlap across
 53 different blocks, heads, and scales. For leveraging this redundancy, following previous work, we
 54 propose sharing the attention output between the conditional and unconditional branches, thereby
 55 skipping the computation for one branch. Combining multidiagonal windowed attention with the
 56 CFG sharing technique, we could reduce 85-90% of attention computation.

57 **Redundancy In Data Precision.** Prior low-bit quantization methods [3, 7] reveal that the high
 58 precision floating-point (FP) representation for neural network weight and activation are redundant.
 59 The Post Training Quantization (PTQ) has proven to be an effective method for both reducing model
 60 size, memory footprint, and computational complexity. Following recent advances in diffusion visual
 61 generation model quantization [18, 19], we apply post training quantization technique to VAR models.
 62 Although W8A8QKV8 quantization achieves satisfying performance. We empirically witness notable
 63 visual quality degradation for lower bit-width (W6A6 and W4A8). Furthermore, we discover that
 64 the quantization is “bottlenecked” by some highly sensitive layers under lower bit-width, and adopt
 65 mixed precision quantization method to preserve these highly sensitive layers at higher bit-width.

66 We summarize the knowledge of our redundancy analysis, the performance-efficiency trade-off, and
 67 deployment feasibility of existing methods in Sec 5.2.

68 2 Attention Redundancy: Multi-Diagonal Window Attention (MDWA)

69 We visualize the attention map for VAR models in Fig. 2-(a). As shown, for most tokens, the attention
 70 map concentrates within local regions at each scale, with more than 80% of the values representing
 71 interactions between spatially distant visual tokens being close to zero. Therefore, replacing the
 72 original global attention with local windowed attention can significantly reduce computation while
 73 preserving the majority of meaningful values in the attention map. Leveraging the unique multi-
 74 diagonal characteristics of the VAR attention map, we propose a specialized multi-diagonal windowed
 75 attention (MDWA) pattern to compress redundancy at the attention map level.

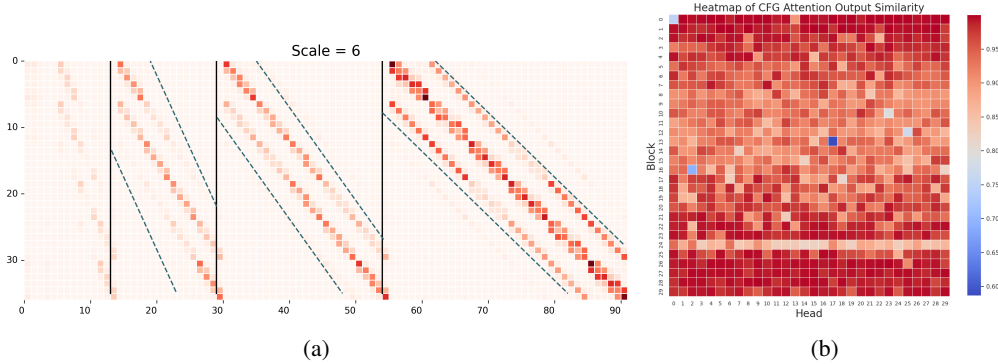


Figure 2: **Attention map characteristics.** (a) **Multi-diagonal concentration.** VAR model’s attention values are concentrated on multiple diagonals, with each diagonal exhibiting a distinct shape across different scales. Consequently, we have designed a separate window attention mechanism for each scale, which we refer to as Multi-Diagonal Window Attention (MDWA). (b) **Similarity** of Attention Outputs between **Conditional** and **Unconditional** Generation.

76 Specifically, considering the VAR model with K scales, each scale containing s_k^2 tokens. For the
 77 k -th scale, the attention mechanism aggregates tokens from the current scale (s_k^2) with all tokens
 78 from previous scale ($\sum_1^k s_i^2$). For example, when $k = 2$, the attention map X has a shape of $[4, 5]$,
 79 where 4 represents the number of visual tokens at the current scale (2^2), and 5 represents the total
 80 tokens from previous scales ($1^2 + 2^2$). As shown in Fig.2-(a), we separate the attention map into
 81 N parts (indicated by vertical black lines) and design a local windowed attention pattern (marked
 82 by blue lines) with window width w . We introduce a metric, R_w , to control the trade-off between
 83 performance and efficiency. The R_w is defined as the division of the summation of all elements
 84 within the window, with respect to the summation of all values in the current part. Since the attention
 85 values are within range $[0, 1]$, the value of R_w could be interpreted as the measurement of “how many
 86 percentage of dominant attention values are contained in the local window”. We gradually increase
 87 the window size from zero until R_w reaches a specific pre-defined ratio R_0 (e.g., 0.95). Table 1
 88 presents the performance efficiency trade-off with different R_0 . When $R_0 = 1$, the attention pattern
 89 falls back to full attention. We further provide the detailed process of the MDWA pattern design.

90 (1) We perform model inference on a subset of the training data and save the attention maps (after
 91 softmax) as a reference for designing the attention pattern.

92 (2) Given an attention map at the k -th scale with the shape $[s_k^2, \sum_1^k s_i^2]$, we first divide it into $k - 2$
 93 parts, where the first part contains $[s_k^2, \sum_1^3 s_i^2]$, and the rest j -th part has the shape $[s_k^2, s_{k-j+1}^2]$. For
 94 each part, we gradually increase the window size w until the ratio R_w reaches a predefined value R_0 .

95 (3) This process is repeated to determine the optimal window size for each scale, block, and head of
 96 the attention map.

97 **Image Evaluation Settings.** We adopt FID [1], IS [10] for fidelity evaluation, and ImageReward [15]
 98 for human preference. Following the original VAR code implementation, we use the 10-scale VAR
 99 with a CFG scale of 4. We generate 8K images on the ImageNet dataset to ensure the stability of the
 100 metric scores.

101 **MDWA implementation details.** In the original VAR design, the s_k values for the 10 scales are
 102 $(1, 2, 3, 4, 5, 6, 8, 10, 13, 16)$. We collect 80 samples in the training set and save their attention maps as
 103 reference. The multi-diagonal windowed attention patterns are designed following the aforementioned
 104 process. Additionally, through analyzing the distribution of attention values, we observe that in the
 105 initial parts of the attention map, certain tokens occasionally exhibit uniformly high attention values
 106 across all tokens. This aligns with the “attention sink” phenomenon described in prior literature [14].
 107 Since the computational cost of these initial parts is relatively low, we retain the full attention pattern
 108 for the first three parts of the attention map.

109 **Experimental Results.** The width of our designed multi-diagonal window attention mechanism was
 110 determined by a threshold setting. We tested different threshold values, including 0.95, 0.9, 0.85,

Table 1: **Performance of MDWA for different Threshold on ImageNet.** Image quality evaluation and Calculation saving for different **Threshold** settings in Multi-Diagonal Window Attention.

| Threshold | FLOPs Saving(%) | FID(↓) | IS(↑) | Image Reward(↑) |
|-----------|-----------------|--------|--------|-----------------|
| 1 | 0.00 | 13.39 | 257.34 | -0.28 |
| 0.95 | 70.34 | 13.47 | 260.95 | -0.28 |
| 0.90 | 73.43 | 13.50 | 261.45 | -0.29 |
| 0.85 | 75.47 | 13.72 | 259.54 | -0.31 |
| 0.80 | 76.82 | 13.77 | 258.45 | -0.34 |
| 0.70 | 79.36 | 13.94 | 254.17 | -0.40 |
| 0.60 | 81.39 | 14.39 | 250.97 | -0.48 |



(a) original (b) MDWA (c) MDWA+ASC

Figure 3: Comparison of original image generation with the techniques of Multi-Diagonal Window Attention (MDWA) and CFG-wise attention sharing (ASC).

111 0.8, 0.7, and 0.6, and evaluated the image quality generated under each threshold. We generated 8k
 112 ImageNet images for evaluation, as shown in Table 1. The threshold of 0.95 yielded the best results,
 113 while a threshold of 0.6 still produced acceptable image quality.

114 3 CFG Redundancy: Attention Sharing across CFG (ASC)

115 Classifier-free guidance (CFG) is widely used for conditional generation [9][8][2], requiring two
 116 model inferences: one with the condition signal and one without. Previous research [17] has explored
 117 reducing the redundancy from the similarity between conditional and unconditional inferences in
 118 diffusion models. Building on this, we investigate similar redundancy in AR-based image generation.
 119 As shown in Fig.2-(b), we observe high similarity between the attention maps of conditional and
 120 unconditional inferences. Based on this, we propose the Attention Sharing across CFG (ASC)
 121 technique, which reuses the attention output from the conditional inference for the unconditional
 122 inference, significantly reducing attention computation cost. Since the vast majority of layers exhibit
 123 high attention map similarity, we reuse the attention maps across the entire network. We will further
 124 explore selectively reusing maps in layers with higher similarity to balance performance and efficiency
 125 in future work.

126 **Experimental Results.** We applied the Attention Sharing across CFG (ASC) technique with the
 127 MDWA technique, the results, as presented in Table 3. The generated images indicate that the loss
 128 introduced by ASC is minimal. In fact, for some metrics, ASC even outperformed the non-shared
 129 attention computation, demonstrating its effectiveness. Combining the MDWA with ASC, we could
 130 achieve 85%-90% attention computation savings with negligible visual quality degradation.

131 4 Data Precision Redundancy: Mixed Precision Quantization

Post Training Quantization (PTQ) has proven to be an efficient and effective model compression
 method [7]. It converts the floating-point data into low-bit integers, the process could be represented
 as:

$$x_q = \text{round}(\text{clamp}((x - z)/s, -2^{B-1}, 2^{B-1}))$$

The s (scale) and z (zero point) are quantization parameters, which are determined offline based on
 stored calibration data with:

$$s = \max(\text{abs}(x))$$

$$z = (\max + \min)/2$$

132 However, we empirically observe that using this straightforward quantization method leads to sig-
 133 nificant quality degradation, even at W8A8QKV8 (weights, activation, and the QKV in attention
 134 are quantized to 8-bit integers). Building on recent advancements in language model quantiza-
 135 tion [13][16], we adopt dynamic quantization parameters for activation quantization, where s and z
 136 are computed online to adapt to diverse activations. Since calculating these quantization parameters
 137 only requires obtaining the maximum and minimum values of the data, the additional computational
 138 cost remains minimal. We apply this dynamic quantization scheme to VAR models, with results
 139 presented in Table 2.

140 While achieving W8A8QKV8 quantization without performance loss, we still observe quality degra-
 141 dation at lower bit-widths (e.g., W4A8QKV8). To investigate the cause, we analyzed the model and
 142 found that quantizing certain layers leads to significant performance drops, while others do not. This
 143 reveals that highly quantization-sensitive layers create a bottleneck for low-bit quantization. As shown
 144 in Fig. 5 in the appendix, our extensive analysis of the VAR model layers indicates that quantizing
 145 the "ffn.fc2" layer to W4A6 causes a disproportionately larger quality degradation compared to
 146 other layers. To address this "bottleneck phenomenon", we propose employing mixed precision
 147 quantization, maintaining higher bit-widths for these particularly sensitive layers.

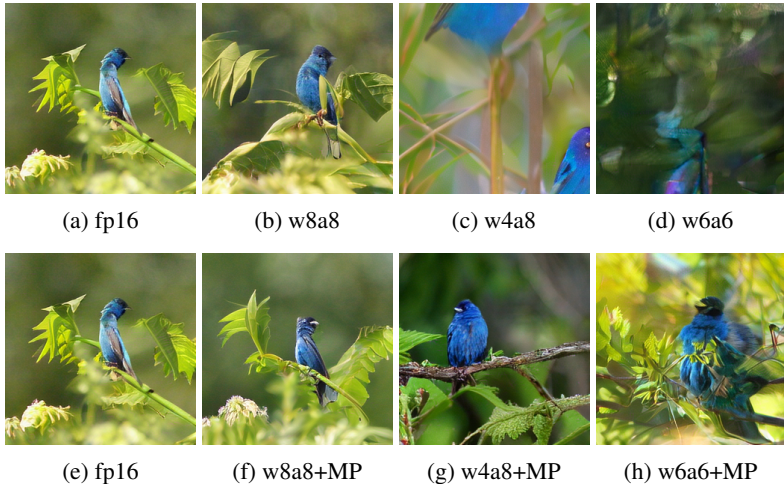


Figure 4: **Comparison of original image, quantized image and quantized image with protection of sensitive layers.** Top row: Naive quantized image exhibit substantial blurring or loss of legible content. Bottom row: A significant improvement in image quality post-quantization.

148 **Quantization Scheme.** We adopt the simple min-max quantization scheme. The quantization
 149 parameters for activation are dynamic and computed online with negligible overhead. The mixed
 150 precision plan are determined offline based on the calibration data.

151 **Experimental Results.** The evaluation scheme are kept consistent with Sec.2. As shown in Table
 152 2 and Fig. 4, both W8A8 and W8A8QKV8 exhibit no performance loss, generating images nearly
 153 identical to those produced with FP16. However, the images generated by W4A8 and W6A6 show
 154 noticeable blurring, underscoring the need for mixed precision quantization. By adopting mixed
 155 precision quantization, both W4A8 and W6A6 experience significant improvements in visual quality
 156 and metric scores. In fact, W4A8 with mixed precision can achieve nearly the same generation quality
 157 as uniform W8A8 quantization.

158 5 Analysis

159 5.1 Ablation Studies

160 As demonstrated in Fig. 3, the introduction of MDWA and ASC results in only a slight performance
 161 degradation (+0.05 FID). Furthermore, replacing the uniform W4A8QKV8 quantization with a
 162 mixed precision scheme significantly reduces performance loss. LiteVAR maintains performance
 163 comparable to the FP16 baseline while effectively compressing redundancy across three dimensions.

Table 2: Performance of image generation on ImageNet under various settings of quantization. Mixed-precision design significantly improves the performance under low bitwidth quantization.

| Bit-width (W/A/QKV) | Mix-Precision (FP16) | FID(↓) | IS(↑) | Image reward(↑) |
|------------------------|----------------------|-----------------------|-------------------------|-----------------------|
| 16/16/16 | – | 13.39 | 257.34 | -0.28 |
| 8/8/8 | – ✓ | 12.71 13.08 | 249.04 253.43 | -0.33 -0.30 |
| 4/8/8 | – ✓ | 54.29 12.82 | 40.71 228.59 | -1.43 -0.41 |
| 6/6/8 | – ✓ | 66.53 18.54 | 26.08 133.13 | -1.68 -0.75 |
| 4/6/8 | – | 111.24 | 9.79 | -2.10 |
| 4/4/8 | – | 133.38 | 6.63 | -2.15 |

Table 3: **Ablation studies of LiteVAR techniques.** When gradually incorporating LiteVAR’s techniques, compressing attention by 85% and reducing the bit width to W4A8QKV8, the generated images are acceptable.

| Method | | | FID | IS | ImageReward |
|--------|-----|----------------|-------|--------|-------------|
| MDWA | ASC | Quant(W/A/QKV) | (↓) | (↑) | (↑) |
| – | – | 16/16/16 | 13.39 | 257.34 | -0.28 |
| ✓ | – | 16/16/16 | 13.47 | 260.95 | -0.28 |
| ✓ | ✓ | 16/16/16 | 13.45 | 248.8 | -0.27 |
| ✓ | ✓ | 4/8/8 | 52.27 | 33.87 | -1.6 |
| ✓ | ✓ | 4/8/8+MP | 13.34 | 224.74 | -0.39 |

164 **5.2 Takeaways for VAR Compression Techniques**

165 **Efficiency Improvement.** The MDWA and CFG-sharing could reduce 85%-90% attention compu-
 166 tation and reduce 80% attention map activation memory cost with negligible computational cost.
 167 Although for current application (ImageNet 256×256), the attention computation and attention
 168 map memory cost is not excessive. However, the attention computation and memory cost grows
 169 quadratically with the token length. For higher resolution (2K) generation, the attention operation
 170 becomes the major bottleneck. In such case, the efficient attention mechanism could significantly
 171 reduce the computation cost (69.6% of the FLOPs), and the memory cost for saving the attention
 172 map (31.07GB). The quantization could effectively reduce both the computational cost and memory
 173 cost of the model. Taking W8A8 as an example, it could reduce 2× of model memory, and achieve
 174 around 1.5× latency speedup.

175 **Efficiency of Compression Methods.** In addition to the efficiency improvement that the compression
 176 method brings, the efficiency of the compression method itself is also critical for practical application.
 177 Therefore, we design **training-free** compression techniques. Unlike many pruning-based methods
 178 that require model fine-tuning, MDWA attention compression eliminates the need for additional
 179 training or large-scale data. Similarly, for post-training quantization, we employ an efficient scheme
 180 that does not rely on gradient-based optimization of quantization parameters.

181 **Deployment Feasibility.** The CFG-sharing technique requires no additional hardware support to
 182 implement, while the MDWA and quantization requires customized CUDA kernels to achieve speedup
 183 and memory savings. For the low-bit quantization, we adopt the commonly used minmax dynamic
 184 quantization scheme, which is supported by many deployment frameworks [4, 20]. The mixed
 185 precision quantization also does not requires additional support other than the W4A8 kernel (which
 186 is also supported by mainstream deployment frameworks).

References

- 187
- 188 [1] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
189 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*
190 *neural information processing systems*, 30, 2017.
- 191 [2] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*
192 *arXiv:2207.12598*, 2022.
- 193 [3] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard,
194 Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for
195 efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer*
196 *vision and pattern recognition*, pages 2704–2713, 2018.
- 197 [4] Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song
198 Han. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving. *arXiv*
199 *preprint arXiv:2405.04532*, 2024.
- 200 [5] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
201 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*
202 *of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- 203 [6] Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar
204 quantization: Vq-vae made simple. *arXiv preprint arXiv:2309.15505*, 2023.
- 205 [7] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen,
206 and Tijmen Blankevoort. A white paper on neural network quantization. *arXiv preprint*
207 *arXiv:2106.08295*, 2021.
- 208 [8] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
209 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3,
210 2022.
- 211 [9] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton,
212 Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.
213 Photorealistic text-to-image diffusion models with deep language understanding. *Advances in*
214 *neural information processing systems*, 35:36479–36494, 2022.
- 215 [10] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.
216 Improved techniques for training gans. *Advances in neural information processing systems*, 29,
217 2016.
- 218 [11] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan.
219 Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint*
220 *arXiv:2406.06525*, 2024.
- 221 [12] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive mod-
222 eling: Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*,
223 2024.
- 224 [13] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han.
225 Smoothquant: Accurate and efficient post-training quantization for large language models.
226 In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.
- 227 [14] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
228 language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- 229 [15] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
230 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
231 *Advances in Neural Information Processing Systems*, 36, 2024.
- 232 [16] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong
233 He. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers.
234 *Advances in Neural Information Processing Systems*, 35:27168–27183, 2022.
- 235 [17] Zhihang Yuan, Pu Lu, Hanling Zhang, Xuefei Ning, Linfeng Zhang, Tianchen Zhao, Shengen
236 Yan, Guohao Dai, and Yu Wang. Ditfastattn: Attention compression for diffusion transformer
237 models. *arXiv preprint arXiv:2406.08552*, 2024.

- 238 [18] Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. Ptq4vit: Post-training
239 quantization for vision transformers with twin uniform quantization. In *European conference*
240 *on computer vision*, pages 191–207. Springer, 2022.
- 241 [19] Tianchen Zhao, Tongcheng Fang, Enshu Liu, Wan Rui, Widyadewi Soedarmadji, Shiyao
242 Li, Zinan Lin, Guohao Dai, Shengen Yan, Huazhong Yang, et al. Vedit-q: Efficient and
243 accurate quantization of diffusion transformers for image and video generation. *arXiv preprint*
244 *arXiv:2406.02540*, 2024.
- 245 [20] Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind
246 Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and
247 accurate llm serving. *Proceedings of Machine Learning and Systems*, 6:196–209, 2024.
- 248 [21] Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenzhe Liu,
249 Lirui Zhao, Fu-Yun Wang, Zhanyu Ma, et al. Lumina-next: Making lumina-t2x stronger and
250 faster with next-dit. *arXiv preprint arXiv:2406.18583*, 2024.

251 **A Appendix / supplemental material**

Visualizing the Sensitivity of various kind of Linear Layers to Bit-Width Reduction.

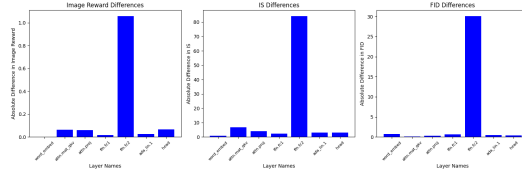


Figure 5: Comparison the impact on image quality of all seven types of linear layers: "word_embed", "attn.mat_qkv", "attn.proj", "ffn.fc1", "ffn.fc2", "ada_lin.1", and "head".

252

253 We observed a particularly noticeable decrease in image quality after quantization for the "ffn.fc2"
 254 layer. To address the quantization bottleneck, we have set the bit width of **ffn.fc2** to **FP16** to safeguard
 255 sensitive layers.

256 **More data for baseline quant:**

Table 4: Performance of image generation on ImageNet under various bitwidths of quantization.

| Bit-width (W/A/QKV) | FID(↓) | | | IS(↑) | | | Image reward(↑) | | |
|------------------------|--------------|--------------|--------------|---------------|---------------|---------------|-----------------|--------------|--------------|
| | Original | Mask | Cfg | Original | Mask | Cfg | Original | Mask | Cfg |
| 16/16/16 | 13.39 | 13.47 | 13.45 | 257.34 | 260.95 | 248.80 | -0.28 | -0.28 | -0.27 |
| 8/8/16 | 12.92 | 13.21 | 13.52 | 252.20 | 258.15 | 244.45 | -0.32 | -0.32 | -0.30 |
| 8/8/8 | 12.71 | 13.02 | 13.38 | 249.04 | 241.02 | 241.04 | -0.33 | -0.37 | -0.29 |
| 4/8/8 | 54.29 | 56.31 | 52.27 | 40.71 | 33.76 | 33.87 | -1.43 | -1.54 | -1.60 |
| 6/6/8 | 66.53 | 68.89 | 62.09 | 26.08 | 24.57 | 28.13 | -1.68 | -1.73 | -1.72 |
| 4/6/8 | 111.24 | 112.79 | 102.44 | 9.79 | 9.52 | 10.48 | -2.10 | -2.13 | -2.10 |
| 4/4/8 | 133.39 | 134.26 | 139.40 | 6.63 | 6.60 | 5.89 | -2.15 | -2.15 | -2.14 |

257 We generated 8,000 images on ImageNet to evaluate the quality of our approach. The bitwidth
 258 designs for the linear layer portion included W16A16 (the original unquantized model), W8A8,
 259 W4A8, W6A6, and W4A6. For the attention computation part, we explored bit-widths of KV8 and
 260 KV16. As shown in the table, quantizing the KV section to a bit-width of 8 has minimal impact on
 261 image quality. When the quantization precision for the linear layer is set to W8A8KV8, the image
 262 quality is comparable to the original floating-point 16-bit (fp16) images. However, W4A8 and W6A6
 263 exhibited significant blurring, and W4A4 resulted in completely illegible images. Subsequently,
 264 we integrated quantization techniques with sparse attention computation to discuss whether the
 265 accuracy could still be maintained. As indicated in the table, the image quality degradation after
 266 sparse computation and ASC (Attention Sharing across CFG) is minimal, demonstrating that we can
 267 significantly reduce computational requirements by approximately 70-90% while ensuring image
 268 quality is preserved.

269 **More data for quantization with mixed-precision design to protect sensitive layers:**

270 Experimental data in 4 reveals that when the weights and activations of linear layers, as well as the
 271 attention computation, are set to 8-bit width, the image quality is essentially preserved. However,
 272 when the weights and activations are designed with lower bit-widths (e.g., W6A6 or W4A8), the
 273 image quality degrades significantly. This is due to the sensitivity of the "ffn.fc2" layer type to
 274 quantization, as illustrated in Figure 5. To address this phenomenon, we set the bit-width of this layer
 275 type to fp16, while maintaining the quantized bit-widths for other layers. We can observe that this
 276 mixed-precision design significantly improves the performance of W6A6 and W4A8, resulting in
 277 noticeably better image quality. For certain metrics (e.g., FID), the W4A8 configuration can even
 278 achieve comparable performance to the baseline W8A8 quantization.

279 **More examples for image generation in different quantization settings.**

Table 5: Performance of image generation on ImageNet under various settings of quantization. Mixed-precision design significantly improves the performance under low bitwidths quantization.

| Bit-width (W/A/QKV) | FID(↓) | | | IS(↑) | | | Image reward(↑) | | |
|------------------------|--------------|--------------|--------------|---------------|---------------|---------------|-----------------|--------------|--------------|
| | Original | Mask | Cfg | Original | Mask | Cfg | Original | Mask | Cfg |
| 16/16/16 | 13.39 | 13.47 | 13.45 | 257.34 | 260.95 | 248.80 | -0.28 | -0.28 | -0.27 |
| 8/8/8 | 12.71 | 13.02 | 13.38 | 249.04 | 241.02 | 241.04 | -0.33 | -0.37 | -0.29 |
| 8/8/8+MP | 13.08 | 13.37 | 13.57 | 253.43 | 251.84 | 244.16 | -0.30 | -0.34 | -0.28 |
| 4/8/8 | 54.29 | 56.31 | 52.27 | 40.71 | 33.76 | 33.87 | -1.43 | -1.54 | -1.60 |
| 4/8/8+MP | 12.82 | 13.23 | 13.34 | 228.59 | 229.58 | 224.74 | -0.41 | -0.44 | -0.39 |
| 6/6/8 | 66.53 | 68.89 | 62.09 | 26.08 | 24.57 | 28.13 | -1.68 | -1.73 | -1.72 |
| 6/6/8+MP | 18.54 | 22.19 | 20.11 | 133.13 | 117.32 | 125.88 | -0.75 | -0.86 | -0.80 |

280 Further examples are presented in the following figures, which compare the original image to both
 281 the quantized image and the quantized image with the enhanced protection of sensitive layers.

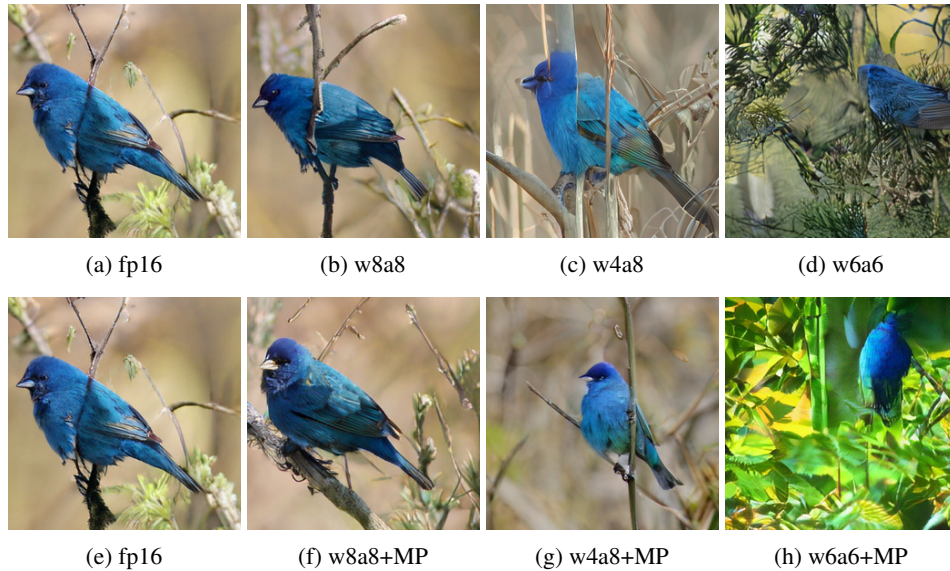


Figure 6: More comparison examples.

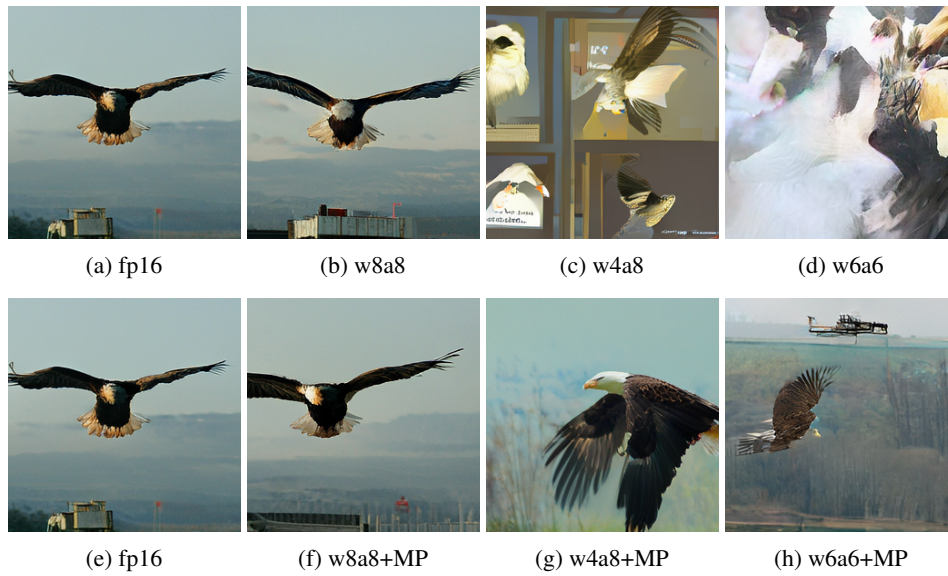


Figure 7: More comparison examples.

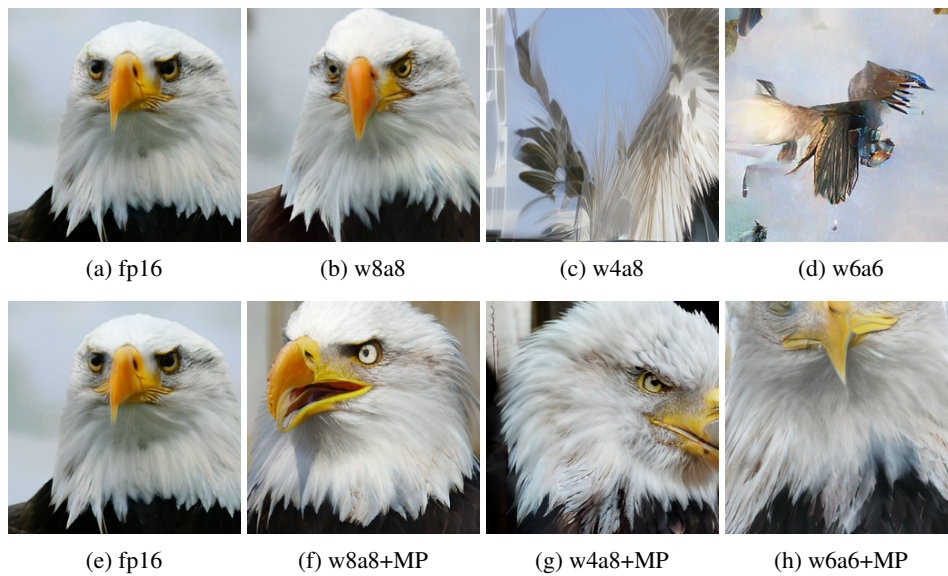


Figure 8: More comparison examples.