

# EVAR: Edge Visual Autoregressive Models via Principled Pruning

Anonymous CVPR submission

Paper ID 5269

001

## Abstract

Recent advances in generative modeling have catalyzed demand for on-device single-image synthesis. However, the stringent compute and memory budgets of resource-constrained edge hardware hinder the deployment of large-scale models. Next-scale visual autoregressive (VAR) models—which predict finer-scale content conditioned on coarser resolutions—offer strong fidelity, generalization, and improved inference efficiency, yet remain costly to run on such devices. We introduce EVAR, an efficient structured-pruning framework tailored to next-scale VAR models and edge deployment. EVAR instantiates a principled pruning paradigm: it couples Optimal Brain Surgeon-guided, Hessian-aware sensitivity estimation with closed-form weight updates, and augments them with scale-aligned calibration and compensation. By grounding pruning decisions in second-order optimality and executing updates analytically, EVAR mitigates compression-induced degradation while preserving next-scale conditioning—turning sparsification from a heuristic into a disciplined procedure. To further address scale-wise gradient and loss imbalance during fine-tuning, we propose Progressive Scale-Aware Distillation (PSAD), which leverages VAR’s multi-scale generative hierarchy to reweight scales and enforce cross-scale consistency in the pruned model. On ImageNet single-image generation benchmarks, EVAR reduces parameter count, memory footprint, and end-to-end latency while retaining competitive generative quality. On an iOS deployment, EVAR further cuts single-image latency from 494 ms to 277 ms (**1.8x speedup**), with FID changing only marginally.

## 1. Introduction

Autoregressive (AR) [7, 21, 42, 51] models have recently made remarkable progress in visual generation, where their strong scalability and generalization enable high-quality image synthesis and editing. However, conventional AR

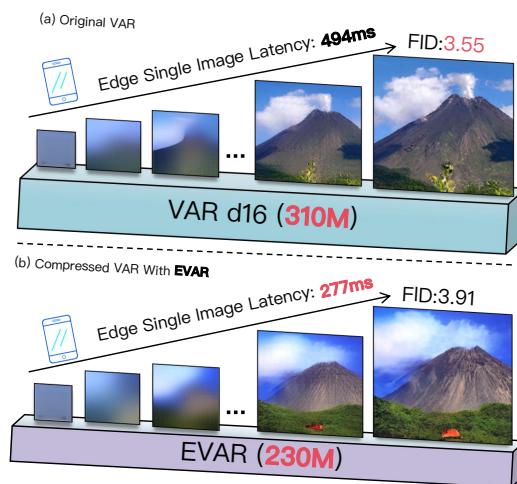


Figure 1. EVAR compresses a next-scale VAR-d16 backbone for edge deployment. (a) Original VAR-d16 with 310M parameters achieves single-image latency of 494 ms with FID 3.55. (b) After EVAR’s structured pruning, the model is reduced to 230M parameters and reaches 277 ms latency on iPad Pro (M4), corresponding to about 1.8× speedup with only ~10% relative FID degradation.

models rely on a next-token prediction paradigm and generate tokens strictly sequentially, so the token-by-token decoding requires numerous steps and leads to significant inference latency.

Visual autoregressive (VAR) modeling [5, 6, 13, 23, 24, 27, 37, 43, 44, 49, 52–54] shifts this paradigm from next-token to next-scale prediction. Instead of emitting tokens one by one, VAR decodes visual content in a coarse-to-fine multi-scale hierarchy, generating many tokens in parallel within each scale. This design drastically reduces the number of decoding steps while retaining the expressive power of AR models, enabling performance comparable to state-of-the-art diffusion models on image synthesis, super-resolution, and inpainting. Yet these models remain large and computationally demanding, making on-device deploy-

037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051

052 ment on resource-constrained edge hardware—where per-  
053 image latency and memory are critical—particularly chal-  
054 lenging. In this work, we focus on single-image generation  
055 in such edge settings.

056 These constraints motivate efficient model compression  
057 techniques [4, 8, 15, 17, 47, 48, 50] that can reduce the  
058 computational and memory footprint without significantly  
059 degrading generative quality. However, existing pruning  
060 methods are mainly developed for classification networks  
061 or large language models and do not account for the multi-  
062 scale structure and generation mechanism of VAR.

063 To address this gap, we propose *EVAR*, a structured com-  
064 pression framework tailored to VAR and edge deployment.  
065 EVAR adopts an OBS-guided [14] adaptive pruning-and-  
066 compensation scheme that assigns block-wise pruning ra-  
067 tios based on Hessian information and explicitly compen-  
068 sates for the effect of removed weights. By coupling pruning  
069 decisions with each block’s internal structure and sen-  
070 sitivity, EVAR substantially reduces parameters and com-  
071 putation while preserving generation quality, making VAR  
072 more suitable for deployment on edge devices.

073 Beyond pruning, we find that next-scale VAR suf-  
074 fers from inherent scale-wise gradient imbalance: high-  
075 resolution scales contain far more tokens than coarse scales,  
076 so standard training objectives over-emphasize fine details  
077 while under-supervising coarse semantics, an issue further  
078 exacerbated after pruning. EVAR addresses this with a  
079 *progressive scale-aware distillation* (PSAD) scheme that  
080 combines a coarse-to-fine distillation curriculum with scale-  
081 aware loss weighting, rebalancing gradients across scales,  
082 and improving the fine-tuning of pruned VAR models.

083 To evaluate EVAR under realistic deployment condi-  
084 tions, we build an end-to-end on-device pipeline for single-  
085 image VAR generation on iOS. We convert pruned models  
086 using Apple’s official toolchain, run them in a Swift-based  
087 iOS application with carefully selected compute units, and  
088 measure end-to-end latency and memory usage. On Ima-  
089 geNet, EVAR delivers competitive generative performance  
090 while enabling efficient mobile and edge deployment of  
091 VAR models.

092 Our contributions can be summarized as follows:

- We introduce *EVAR*, the first OBS-guided adaptive pruning and compensation framework tailored for VAR models, providing a principled structured pruning pipeline for edge deployment.
- We propose progressive scale-aware distillation (PSAD), which addresses scale-wise gradient imbalance in next-scale VAR and improves the fine-tuning of pruned models via a coarse-to-fine, scale-weighted distillation.
- Empirically, we build a practical on-device deployment and evaluation pipeline for single-image VAR generation on iOS with CoreML inference engine, reporting 1.8× speedup with comparable generation quality.

## 2. Related Work

### 2.1. Efficient Auto-regressive Image Generation

Autoregressive (AR) models decompose the joint distribution of an image into a product of conditional distributions and have been a long-standing backbone of generative modeling. Early work, such as PixelRNN and Pixel-CNN [31, 45] generates pixels sequentially with recurrent or convolutional architectures, but suffers from slow sampling. With the advent of Transformer-based language models [2, 32, 33, 35, 36, 46], large-scale GPT-style decoders have been adopted for images, enabling global context modeling and strong scalability. Recent systems, including VQGAN and RQ-Transformer [7, 19], diffusion–AR hybrids [11, 55], and GPT-like image generators such as LlamaGen and Lumina-mGPT [26, 42] further push fidelity by operating in discrete token spaces with powerful next-token prediction.

To reduce the sequential bottleneck, a line of work studies partially parallel or masked decoding, e.g., MaskGIT and MAR [3, 22], and scalable diffusion backbones such as DiT and LDM [34, 38]. CoDe [6] boosts efficiency with minimal quality loss via collaboration: large drafter for small-scale low-frequency, small refiner for large-scale high-frequency details. This design achieves competitive quality with substantially fewer decoding steps, and has been extended to strong AR baselines such as PAR and AiM [21, 51]. However, most of these works target fidelity and scalability on powerful GPUs; compression and practical deployment of next-scale VAR on resource-constrained edge devices remain largely unexplored, which is the focus of EVAR.

### 2.2. Network Pruning

Network pruning [12, 30] is a standard way to reduce computation and memory by removing redundant parameters from a pre-trained network. According to granularity, methods are typically divided into *unstructured* and *structured* pruning. Unstructured pruning [9, 10, 14, 18, 39, 40] zeros out individual weights according to saliency scores, enabling fine-grained parameter removal at the level of each matrix entry. This granularity often yields lower accuracy degradation at a given sparsity compared to structured pruning. However, the resulting irregular sparsity is difficult to exploit on general-purpose hardware due to poor alignment with SIMD/SIMT execution, irregular memory access patterns, and limited kernel support; consequently, end-to-end speedups are frequently modest despite high nominal sparsity. Structured pruning [1, 8, 16, 20, 25, 28, 41], removes entire channels or heads, yielding dense, smaller subnetworks that better align with standard conv/GEMM kernels and thus more directly translate into wall-clock acceleration.

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

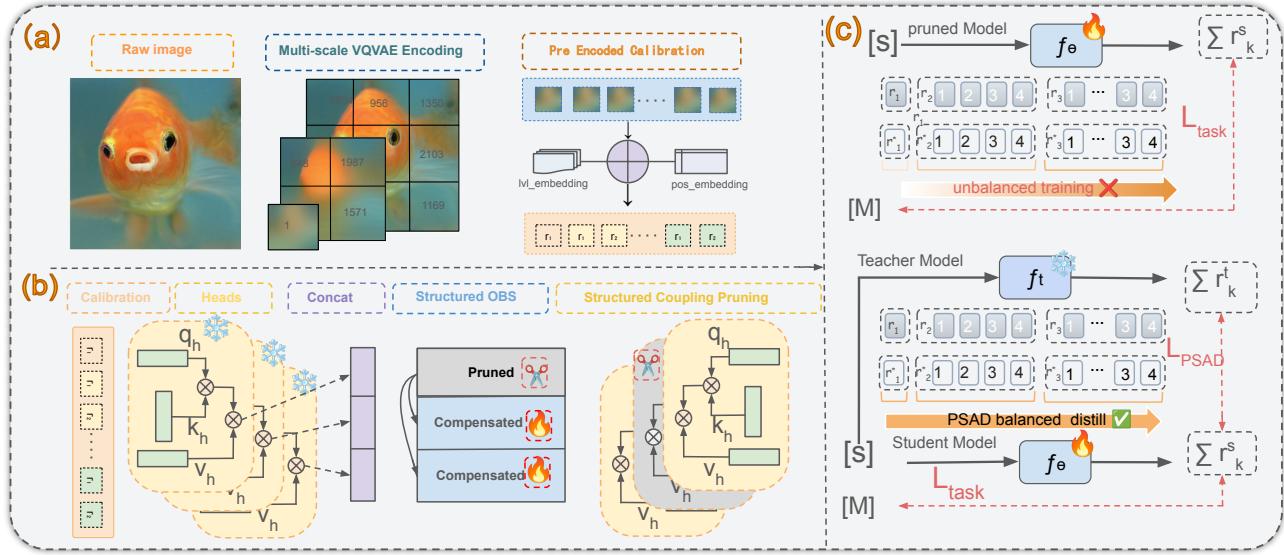


Figure 2. (a) *Pre-encode calibration*: a raw image is encoded by the VQ-VAE into a multi-scale residual pyramid, which is then converted into next-scale VAR tokens with level and positional embeddings to form the pre-encoded calibration set. (b) *OBS-guided structured pruning*: using the input covariance of the output projection  $O$  to approximate its Hessian, EVAR applies structured OBS to prune attention heads, compensates the remaining heads in closed form, and finally removes the corresponding coupled heads in the shared Q/K/V projections. (c) *Progressive Scale-Aware Distillation (PSAD)*: standard fine-tuning or vanilla distillation suffers from scale-wise loss imbalance in next-scale VAR, whereas PSAD reweights and schedules the multi-scale distillation loss to balance supervision across scales, yielding state-of-the-art post-pruning recovery.

156

### 3. Method

157

#### 3.1. Preliminaries on VAR

158

Visual Autoregressive (VAR) models [44] reformulate autoregressive image modeling by shifting from traditional next-token prediction to a next-scale prediction scheme. Instead of generating an image token by token, VAR divides the input feature map  $f \in \mathbb{R}^{h \times w \times C}$  into  $K$  token maps  $(r_1, r_2, \dots, r_K)$  at multiple resolutions, where the spatial resolution increases with the scale index and the final token map  $r_K$  recovers the original feature-map resolution.

166

The joint distribution over the multi-scale token maps is factorized as

168

$$p(r_1, r_2, \dots, r_K) = \prod_{k=1}^K p(r_k | r_1, \dots, r_{k-1}), \quad (1)$$

169

where  $r_k \in [V]^{h_k \times w_k}$  denotes the token map at scale  $k$ , and  $(r_1, \dots, r_{k-1})$  provides the coarse-to-fine context for predicting  $r_k$ . At each autoregressive step  $k$ , all tokens in  $r_k$  are generated in parallel, conditioned on the previous scales and their positional embeddings.

174

#### 3.2. OBS-Guided Structured Pruning for VAR

175

Network pruning reduces model size and inference cost by removing redundant or less important parameters. We focus

on post-training *structured* pruning, where entire channels or heads are removed so that the resulting model remains compatible with hardware acceleration. The central challenge is to reduce parameters while controlling the degradation in generation quality, especially for next-scale VAR models with long token sequences.

177

178

179

180

181

182

**Layer-wise OBS formulation.** Following Optimal Brain Surgeon (OBS), we decompose the global pruning problem into layer-wise subproblems and minimize the discrepancy between the original and pruned layer outputs in a local  $\ell_2$  sense. For the  $l$ -th layer with weight matrix  $W_l$  and input activation  $X_l$ , we consider the objective:

$$\min_{\hat{W}_l} \|W_l X_l - \hat{W}_l X_l\|_2^2, \quad (2)$$

183

184

185

186

187

188

subject to a target pruning ratio on  $W_l$ . A second-order Taylor expansion around  $W_l$  yields a quadratic approximation of the loss in terms of the weight perturbation  $\Delta W_l$ , with layer-wise Hessian:

$$H_l = 2X_l X_l^\top. \quad (3)$$

194

For a scalar layer  $W_l$ , OBS iteratively identifies and removes the weight  $w_p^*$  that has the smallest effect on the objective output, and applies a closed-form compensation to

195

196

197

198

update remaining weights:

$$w_p^* = \arg \min_{w_p} w_p^2 (H_l^{-1})_{p,p}, \quad \delta_p = -\frac{w_p^*}{(H_l^{-1})_{p,p}} (H_l^{-1})_{:,p}, \quad (4)$$

where  $(H_l^{-1})_{p,p}$  and  $(H_l^{-1})_{:,p}$  denote the  $p$ -th diagonal entry and column of  $H_l^{-1}$ , respectively. After pruning  $w_p$  to zero,  $\delta_p$  is added to the remaining weights.

**ExactOBS: row-wise formulation and inverse updates.** To make OBS practical, we follow ExactOBS [9] and exploit two structural properties of the layer-wise objective. Dropping the layer index  $l$  for clarity, let  $W \in \mathbb{R}^{d_{\text{row}} \times d_{\text{col}}}$  and  $Y = WX$  be the dense layer output. The objective in Eq. (2) can be rewritten as a sum of row-wise errors:

$$\|WX - \hat{W}X\|_2^2 = \sum_{i=1}^{d_{\text{row}}} \|W_{i,:}X - \hat{W}_{i,:}X\|_2^2. \quad (5)$$

Removing a scalar entry  $[W]_{i,j}$  only affects the error of the  $i$ -th output row  $Y_{i,:}$ ; there is no Hessian interaction between different rows. Thus, each row can be treated as an independent least-squares problem with a  $d_{\text{col}} \times d_{\text{col}}$  Hessian

$$H = 2XX^\top, \quad (6)$$

which is shared across all rows. This observation reduces the problem from a  $d \times d$  Hessian to a much smaller  $d_{\text{col}} \times d_{\text{col}}$  matrix.

The second ingredient of ExactOBS is an efficient update of the inverse Hessian when a single parameter  $p$  is removed. Let  $H$  be invertible with inverse  $H^{-1}$ , and let  $H_{-p}$  denote the principal submatrix obtained by removing row and column  $p$  from  $H$ . As shown in [9], the inverse  $H_{-p}^{-1}$  can be obtained directly from  $H^{-1}$  via a single Gaussian-elimination-style update:

$$H_{-p}^{-1} = (H^{-1} - \frac{1}{(H^{-1})_{p,p}} H_{:,p}^{-1} (H^{-1})_{p,:})_{-p}, \quad (7)$$

where  $(\cdot)_{-p}$  denotes removing row and column  $p$ . Intuitively, we eliminate the influence of parameter  $p$  from the inverse and then drop the corresponding row and column. This update costs  $O(d_{\text{col}}^2)$  time and does not require recomputing any matrix inverses.

**Structured OBS for VAR blocks.** For next-scale VAR, we extend the above OBS formulation to operate on structured units in the pre-encode Transformer blocks at the last scale. Let  $W^{(S)}$  denote a weight matrix in such a block and  $H^{(S)}$  its Hessian. We adopt column pruning as the basic operation and treat groups of columns as structured units (e.g., attention heads or FFN channels). For a column index  $p$ , the structured OBS criterion and compensation can be written as

$$W_{:,p}^{(S)} = \arg \min_{W_{:,p}^{(S)}} \|W_{:,p}^{(S)}\|_2^2 (H^{(S)})_{p,p}^{-1}, \quad (8)$$

$$\Delta^{(S)} = -\frac{W_{:,p}^{(S)}}{(H^{(S)})_{p,p}^{-1}} (H^{(S)})_{p,:}^{-1}, \quad (9)$$

where  $(H^{(S)})_{p,:}^{-1}$  is the  $p$ -th row of  $(H^{(S)})^{-1}$  and  $\Delta^{(S)}$  is a compensation matrix with the same shape as  $W^{(S)}$ . Following common practice, we prune attention blocks and FFNs as basic units: pruning columns in the attention output projection and FFN down-projection at the last scale reduces the number of heads and intermediate channels, thereby shrinking the model size.

Directly applying stepwise column pruning with OBS is still expensive for VAR due to the large intermediate dimensions and the interdependence between attention heads. We therefore adopt an *iterative unit-wise pruning-and-compensation scheme*: (i) define attention heads and FFN channels as structured units; (ii) estimate the OBS error for each unit using the Hessian statistics from the pre-encode calibration set (Sec. 3.3); (iii) prune the least important unit and apply the compensation update in Eq. (9); and (iv) repeat until the target sparsity is reached. For FFN layers, we further use a dynamic grouped strategy that prunes small groups of low-score channels at a time, starting with larger groups and gradually decreasing the group size as pruning progresses. This improves efficiency while keeping the solution close to the column-wise optimum. Empirically, this structured OBS scheme provides strong compression for VAR while maintaining high image generation quality.

### 3.3. Pre-Encode Calibration for Next-Scale VAR

OBS-guided pruning in a post-training setting relies critically on a calibration set: layer-wise Hessian estimates  $H_l$  are constructed from the activations induced by the calibration samples, and inaccurate calibration directly translates into suboptimal sensitivity estimates and poor training-free performance. For next-scale VAR, naive choices such as using tokens sampled from the model’s own autoregressive generation pipeline (conditioned on prompts or partial inputs) lead to a severe mismatch between the calibration distribution and the real deployment regime.

To address this, we propose a *pre-encode calibration* strategy tailored to next-scale VAR. Instead of calibrating on model-generated tokens, we construct calibration samples by encoding real images through the exact VAE and residual pyramid used by VAR at training and inference time. Concretely, given a real image  $x$ , we first obtain a latent feature map  $f \in \mathbb{R}^{h \times w \times C}$  via the VAE encoder. We then apply the VAR residual pyramid to decompose  $f$  into multi-scale residual feature maps, and quantize each scale using VAR’s codebook and interpolation–lookup procedure to obtain the discrete token maps  $(r_1, \dots, r_K)$ . Finally, we attach the same positional and scale embeddings as in

291 VAR’s training pipeline, yielding complete multi-scale token inputs that faithfully reflect the inference-time distribution.  
 292  
 293

294 We use this pre-encode calibration set to estimate Hessian  
 295 statistics for the Transformer blocks at the last scale  
 296 and to drive the OBS-based structured pruning described in  
 297 Sec. 3.2. Because the calibration tokens are derived from  
 298 real images and respect VAR’s multi-scale residual hierar-  
 299 chy, the resulting Hessian estimates are substantially more  
 300 stable and informative for sensitivity analysis.

### 301 3.4. Progressive Scale-Aware Distillation (PSAD)

302 Next-scale VAR inherently exhibits *scale-wise gradient im-  
 303 balance*: Higher-resolution scales contain far more to-  
 304 kens than coarse scales, so unweighted objectives over-  
 305 emphasize fine scales and under-supervise coarse seman-  
 306 tics. Pruning reduces model capacity and amplifies this  
 307 bias, making global structure harder to preserve for the  
 308 compressed model. We address this issue with *Progres-  
 309 sive Scale-Aware Distillation* (PSAD), which couples a coarse-  
 310 to-fine curriculum with scale-aware weighting to rebalance  
 311 gradients across scales during post-pruning adaptation.

312 **Setup.** Let  $K$  be the number of scales and  $L$  the full se-  
 313 quence length after concatenating tokens from all scales.  
 314 Scales are indexed by  $i \in \{0, \dots, K-1\}$  with token  
 315 boundaries  $\mathbf{b} = [b_0, \dots, b_K]$  and per-scale counts  $\mathbf{n} =$   
 316  $[n_0, \dots, n_{K-1}]$ , where we set

$$317 b_0 = 0, \quad b_K = L, \quad n_i = b_{i+1} - b_i, \quad \sum_{i=0}^{K-1} n_i = L. \quad (10)$$

318 Thus scale  $i$  occupies positions  $\ell \in (b_i, b_{i+1}]$  in the flat-  
 319 tened sequence. Let  $B$  denote the batch size. We use the  
 320 original (unpruned) VAR as a teacher and the pruned VAR  
 321 as a student. With temperature  $\tau > 0$ , teacher and student  
 322 distributions at batch index  $b$  and position  $\ell$  are

$$323 p_T^{(b, \ell)} = \text{softmax}(\mathbf{z}_T^{(b)}[\ell, :] / \tau), \quad p_S^{(b, \ell)} = \text{softmax}(\mathbf{z}_S^{(b)}[\ell, :] / \tau),$$

324 where  $\mathbf{z}_T^{(b)}, \mathbf{z}_S^{(b)}$  denote teacher and student logits, respec-  
 325 tively. The scale-level KL divergence (normalized by  $n_i$ ) is  
 326 defined as

$$327 \mathcal{L}_{\text{KL}}^{(i)} = \frac{1}{B n_i} \sum_{b=1}^B \sum_{\ell=b_i+1}^{b_{i+1}} \text{KL}(p_T^{(b, \ell)} \| p_S^{(b, \ell)}). \quad (11)$$

328 We include a factor  $\tau^2$  in the distillation loss below to com-  
 329 pensate for the  $1/\tau^2$  scaling of gradients with temperature.

330 **1) Discrete stage-wise progression.** We first describe  
 331 a discrete coarse-to-fine schedule that progressively “un-  
 332 locks” scales over training. Define the highest unlocked

333 scale at optimization step  $t$  as

$$334 s(t) = \max\{i \in \{0, \dots, K-1\} \mid t \geq t_i\}, \quad (12)$$

334  
 with  $0 = t_0 < t_1 < \dots < t_{K-1} < T$ ,

335 where  $t_i$  are pre-defined milestones and  $T$  is the total num-  
 336 ber of steps. Let  $\gamma_i(t) = \mathbb{1}[i \leq s(t)]$  indicate whether  
 337 scale  $i$  is active at step  $t$ . The progressive, scale-weighted  
 338 distillation loss is then

$$339 \mathcal{L}_{\text{PSAD}}(t) = \tau^2 \sum_{i=0}^{K-1} \gamma_i(t) w_i \mathcal{L}_{\text{KL}}^{(i)}, \quad (13)$$

340 where  $w_i$  are per-scale weights (specified below). The task  
 341 loss aggregates only the unlocked tokens:

$$342 \mathcal{L}_{\text{task}}(t) = -\frac{1}{B b_{s(t)+1}} \sum_{b=1}^B \sum_{\ell=1}^{b_{s(t)+1}} \log p_S(x_\ell^{(b)} \mid x_{<\ell}^{(b)}, y^{(b)}), \quad (14)$$

343 where  $x_\ell^{(b)}$  denotes the ground-truth token at position  $\ell$  in  
 344 sample  $b$ , and  $y^{(b)}$  denotes the conditioning (e.g., class la-  
 345 bel). The total objective is

$$346 \mathcal{L}_{\text{total}}(t) = \alpha \mathcal{L}_{\text{task}}(t) + \beta \mathcal{L}_{\text{PSAD}}(t), \quad (15)$$

347 with  $\alpha$  and  $\beta$  balancing data loss and distillation.

348 **2) Continuous soft progression.** The hard indicator  $\gamma_i(t)$   
 349 can lead to abrupt changes in the loss when a new scale  
 350 is unlocked. To smooth these transitions, we replace  $\gamma_i(t)$   
 351 with a continuous gate

$$352 \tilde{\gamma}_i(t) = \min\left\{1, \max\left\{0, \frac{t - t_i}{\Delta_i}\right\}\right\}, \quad (16)$$

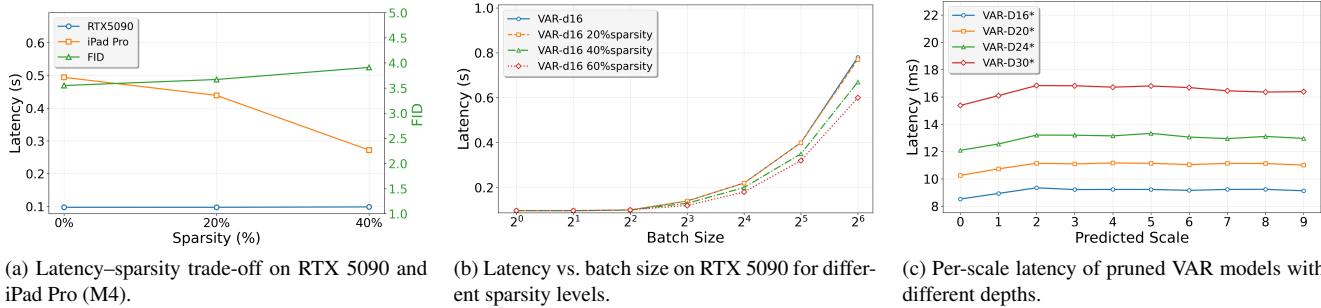
353 where  $\Delta_i > 0$  controls the ramp-up duration for scale  $i$ .  
 354 This yields a soft version of PSAD:

$$355 \mathcal{L}_{\text{PSAD}}^{\text{soft}}(t) = \tau^2 \sum_{i=0}^{K-1} \tilde{\gamma}_i(t) w_i \mathcal{L}_{\text{KL}}^{(i)}. \quad (17)$$

356 In practice, we find that the soft schedule improves stabil-  
 357 ity over the discrete variant, especially when the model is  
 358 heavily pruned.

359 **3) Progressive weights and gradient balancing.** Token-  
 360 length imbalance implies that, under a naïve uniform  
 361 weighting, high-resolution scales with large  $n_i$  dominate  
 362 the gradients. To counter this, we use monotonically de-  
 363 creasing scale weights, e.g.,

$$364 w_i = 2.0 - 0.2i, \quad i = 0, \dots, K-1, \quad (18)$$



(a) Latency-sparsity trade-off on RTX 5090 and iPad Pro (M4).

(b) Latency vs. batch size on RTX 5090 for different sparsity levels.

(c) Per-scale latency of pruned VAR models with different depths.

Figure 3. Latency behavior of EVAR under different sparsity and deployment settings. (a) As sparsity increases, single-image latency on the RTX 5090 remains almost unchanged, whereas latency on iPad Pro (M4) drops substantially while FID increases only mildly, indicating that EVAR makes single-image inference increasingly compute-bound on iOS with acceptable quality loss. (b) On the RTX 5090, dense and pruned VAR-d16 models show nearly identical latency at batch size 1, and only diverge at larger batch sizes, suggesting that single-image inference on the GPU is predominantly memory-bound. (c) For different VAR depths (VAR-D16\*, VAR-D20\*, VAR-D24\*, VAR-D30\*), per-scale latency for single-image inference is nearly flat across predicted scales, and pruning (indicated by \*) does not introduce scale-dependent latency spikes, consistent with our width-oriented design.

which is consistent with the heuristic  $w_i \propto 1/n_i$ . In our experiments with  $K = 10$ , this yields positive and monotonically decreasing weights. The effective per-scale gradient magnitude scales roughly as

$$g_i(t) \propto (\gamma_i(t) \text{ or } \tilde{\gamma}_i(t)) \cdot w_i \cdot n_i \cdot \|\nabla_{\theta_S} \mathcal{L}_{KL}^{(i)}\|.$$

The combination of progressive gates  $(\gamma_i, \tilde{\gamma}_i)$  and decreasing  $w_i$  substantially reduces the dominance of high-resolution scales in the gradient, leading to more balanced multi-scale updates. In combination with the coarse-to-fine schedule, PSAD thus rebalances training signals across scales and significantly improves reconstruction fidelity and overall generative quality compared with conventional fine-tuning, particularly in the post-pruning setting.

## 4. Experimental Results

### 4.1. Experimental Setups

We evaluate EVAR on the ImageNet-1K [29] dataset using images of resolution  $256 \times 256$ , following exactly the same data preprocessing and augmentation pipeline as the original VAR model [44]. For a fair comparison, we adopt the VAR\_d16 model as the base backbone on the ImageNet  $256 \times 256$  conditional generation benchmark, and compare against state-of-the-art image generation model families.

We use class-balanced sampling with one image per class as the pre-encoded calibration set. Enlarging the calibration set yields no significant pruning gains, so we adopt “one per class” as the default.

After pruning, all models are fine-tuned for 40 epochs using AdamW with the same optimizer hyper-parameters, learning-rate schedule, and learning-rate / weight-decay annealing as in the original VAR training. Fine-tuning is performed on 8 NVIDIA 5090 GPUs. Finetuning the pruned model for one epoch takes slightly over one hour.

For edge deployment, the pruned models are further converted and deployed on an iPad Pro (M4), where we measure on-device latency for single-image generation.

### 4.2. Main Results

We compare our method with state-of-the-art pruning and compression approaches. Table 1 reports model size, parameter count, FLOPs, and storage reduction before and after pruning. Within the VAR family, we take VAR-d16 as our base model and report results for two pruned variants obtained with EVAR at 20% and 40% structured sparsity. At 20% pruning (EVAR, 270M parameters), our method maintains a FID of 3.67 compared to 3.55 for the dense VAR-d16, while reducing parameters from 310M / 10.97 GB to 270M. At 40% pruning (EVAR, 230M parameters), EVAR further compresses the model to 230M parameters with a FID of 3.91, which remains competitive given the substantial reduction in model size and decoding cost.

As summarized in Tab. 2, EVAR (ours) consistently outperforms a range of structured pruning baselines across sparsity levels and fine-tuning regimes. Under the most challenging setting of 40% sparsity without any fine-tuning, all methods suffer large quality drops, but EVAR still achieves the best FID/IS (64.19 / 19.52), whereas other methods degrade much more severely (e.g., FID > 140 for OBA, Taylor, and LLM-Pruner). At 20% sparsity in the training-free regime, EVAR attains a substantially lower FID of 8.86 with strong precision/recall, while all baselines remain in a much worse range (FID 27.30–152.29), indicating that OBS-guided pruning with pre-encode calibration is particularly effective for moderate compression.

The row “Inference set” isolates the effect of calibration: it uses the same OBS-based pruning as EVAR but replaces our pre-encode calibration with tokens generated by VAR’s own inference. Its performance is consistently worse than

Table 1. Generative performance on class-conditional ImageNet-256. “Steps” means the number of model inference to generate one image.

Type	Model	Parameters	Pruning Rate	Steps	FID↓	IS↑	Precision↑	Recall↑
VAR	VAR-d16 [44]	310M	–	10	3.55	274.4	0.84	0.51
	VAR-d20 [44]	600M	–	10	2.95	302.6	0.83	0.56
	VAR-d24 [44]	1.0B	–	10	2.33	312.9	0.82	0.59
AR	VQGAN-re [7]	1.4B	–	256	5.20	280.3	–	–
	RQ-Trans.-re [19]	3.8B	–	64	3.80	323.7	–	–
	LlamaGen-L [42]	343M	–	576	3.07	256.1	0.83	0.52
	LlamaGen-XL [42]	775M	–	576	2.62	244.1	0.80	0.57
	LlamaGen-XXL [42]	1.4B	–	576	2.62	244.1	0.80	0.57
	PAR-L [51]	343M	–	147	3.76	218.9	0.81	0.60
pruned VAR-d16	PAR-XL [51]	775M	–	147	2.61	259.2	0.80	0.62
	PAR-XXL [51]	1.4B	–	147	2.35	263.2	0.80	0.62
	AiM-L [21]	350M	–	256	2.83	244.6	0.82	0.55
	AiM-XL [21]	763M	–	256	2.56	257.2	0.82	0.57
	LLM-pruner [28]	230M	40%	10	4.21	53.92	0.81	0.50
	OBA [41]	230M	40%	10	4.19	53.43	0.83	0.47
pruned VAR-d16	EVAR (ours)	270M	20%	10	3.67	57.78	0.81	0.51
	EVAR (ours)	230M	40%	10	3.91	57.23	0.81	0.51

Table 2. Comparison of pruning methods under different sparsity and fine-tuning regimes.

Method	40% sparsity, training-free				20% sparsity, training-free				40% sparsity, 1-epoch			
	FID	IS	Precision	Recall	FID	IS	Precision	Recall	FID	IS	Precision	Recall
EVAR (ours)	64.19	19.52	0.35	0.53	8.86	48.41	0.74	0.51	4.86	55.61	0.82	0.47
Inference set	138.91	6.39	0.09	0.07	56.87	19.35	0.36	0.57	9.62	44.91	0.77	0.46
LLM-pruner [28]	153.82	5.62	0.09	0.02	35.21	28.95	0.41	0.51	10.82	41.35	0.76	0.45
OBA [41]	142.82	6.35	0.12	0.04	30.54	29.48	0.45	0.57	9.85	44.35	0.78	0.46
Magnitude [12]	180.82	3.90	0.17	0.01	152.29	5.77	0.10	0.01	12.69	25.84	0.65	0.38
Taylor [30]	145.66	6.78	0.12	0.05	27.30	31.48	0.48	0.59	8.35	46.87	0.77	0.46

431 EVAR (e.g., FID 56.87 vs. 8.86 at 20% sparsity), showing  
432 that structured OBS compensation only works well when  
433 the Hessian is estimated from a precise dataset-aligned cali-  
434 bration set. After one epoch of post-pruning training, EVAR  
435 further improves to FID 4.86 at 40% sparsity, outperforming  
436 all baselines (best non-EVAR FID 8.35 for Taylor), and con-  
437 firming that combining pre-encode calibration, OBS-guided  
438 structured pruning, and our distillation scheme yields the  
439 strongest overall trade-off between compression and gener-  
440 ative quality.

441 **Recent top-performing methods: LLM-Pruner and**  
442 **OBA.** Since existing pruning methods for Transformers  
443 and large language models do not report results on VAR,  
444 we reimplement two structured pruning baselines on VAR-  
445 d16 using the public code: LLM-Pruner [28] and OBA [41].  
446 For both methods, we prune VAR-d16 to a comparable 40%  
447 structured sparsity level and fine-tune for the same number

of epochs as EVAR to ensure a fair comparison.

448 LLM-Pruner is applied with its default `param_mix`  
449 salience metric, which combines gradient and accumulated-  
450 gradient information to rank structured units. OBA is used  
451 in its structured setting, where sensitivity is estimated via  
452 Hessian–vector products over inter-layer connectivity and  
453 used to prune attention heads and FFN channels in VAR-  
454 d16. The corresponding results for LLM-Pruner, OBA, and  
455 our EVAR models are reported in the last block of Table 1.

### 4.3. Ablation Study

458 We ablate the distillation components in EVAR on  
459 ImageNet-1K single-image generation to quantify their ef-  
460 fect on post-pruning recovery. Table 4 reports FID, iOS  
461 single-image latency, and parameter count for the dense  
462 VAR-d16 baseline, the training-free pruned model, and  
463 successive additions of fine-tuning and distillation compo-  
464 nents. Training-free pruning reduces latency from 494 ms

Device	Model	Compute units	Latency (ms)	#CPU ops	#GPU ops	#NPU ops
iPad Pro (M4)	EVAR	CPU+GPU	277 ± 2	10 (random_category)	8818	0
iPad Pro (M4)	EVAR	CPU+NPU	1090 ± 2	274	0	8554
iPad Pro (M4)	EVAR	All	465 ± 2	129	3103	5596

Table 3. Operator-level unit assignment in Core ML for a pruned VAR backbone (EVAR) on iPad Pro (M4). Frequent switches across CPU/GPU/NPU due to NE-unsupported ops inflate latency; pinning to CPU+GPU is fastest for single-image inference.

Method	FID ↓	iOS (ms) ↓	Params (M)
Dense	3.55	494	310
Pruned (training-free)	64.00	277	230
+ Fine-tuning (FT)	4.25	277	230
+ Vanilla KD	4.26	277	230
+ Scale-aware weights	3.97	277	230
+ Progressive (PSAD, ours)	<b>3.91</b>	277	230

Table 4. Ablation of the distillation components in EVAR. All post-pruning variants share the same pruned VAR-d16 backbone, hence identical latency and parameter counts under the same iOS deployment setting.

to 277 ms and parameters from 310M to 230M, but severely degrades FID (64.0). Plain fine-tuning and vanilla KD recover most of the quality ( $\text{FID} \approx 4.25$ ) while keeping the same latency/parameter budget, yet still underperform the dense model. Adding *scale-aware weights* further improves FID to 3.97, and enabling the full *Progressive Scale-Aware Distillation* (PSAD) yields the best result (FID 3.91), nearly closing the gap to the dense baseline under the same 277 ms / 230M on-device configuration.

#### 4.4. Edge Deployment

**Core ML conversion and runtime.** For on-device evaluation, we convert the pruned VAR-d16 models from PyTorch to Core ML (.mlmodel) using `coremltools`, and run them on Apple silicon in a Swift-based single-image generation app on an iPad Pro (M4). The Core ML runtime performs per-operator scheduling and can *automatically switch compute units*—CPU, GPU, and Neural Engine (NE/NPU)—whenever an operator is unsupported on the currently selected unit.

**Compute-unit behavior.** Next-scale VAR contains several operators (e.g., positional embeddings and sampling-related ops) that are not supported on the NE, so enabling the NE causes frequent fallbacks and handoffs between CPU, GPU, and NE at batch size 1. Table 3 reports an operator-level breakdown for EVAR on an iPad Pro (M4) under three Core ML settings: CPU+GPU, CPU+NPU, and All. Although CPU+NPU assigns most operators to the NE, it incurs many CPU ops (274) and yields the highest latency ( $1090 \pm 2$  ms). The All setting mixes all three units (129 CPU ops, 3103 GPU ops, 5596 NPU ops) and is also slower ( $494 \pm 2$  ms) due to frequent cross-unit trans-

sitions. In contrast, pinning execution to CPU+GPU avoids NE-unsupported operators entirely, yields a simple two-unit schedule (10 CPU ops, 8818 GPU ops), and achieves the lowest latency ( $277 \pm 2$  ms) for single-image inference.

**Latency analysis.** Figure 3 summarizes the latency behavior of EVAR. In Fig. 3(a), increasing sparsity barely changes single-image latency on the RTX 5090, but significantly reduces latency on iPad Pro (M4) with only a mild FID increase, showing that EVAR is effective in the edge compute-bound regime. Fig. 3(b) shows that on the RTX 5090, dense and pruned VAR-d16 models have almost identical latency at batch size 1 and diverge only at larger batch sizes, indicating that single-image inference on GPU is largely memory-bound. Fig. 3(c) reports per-scale latency for different VAR depths (VAR-D16\*, VAR-D20\*, VAR-D24\*, VAR-D30\*); the curves are nearly flat across predicted scales and pruning (\*) does not introduce scale-dependent latency spikes, consistent with our width-oriented design.

#### 5. Conclusion

This paper introduces *EVAR*, a principled structured pruning framework for next-scale visual autoregressive (VAR) models. EVAR combines a pre-encode calibration pipeline with OBS-guided head/channel pruning and closed-form compensation, enabling training-free compression that respects the multi-scale conditioning structure of VAR. To further recover generative quality after pruning, we introduced *Progressive Scale-Aware Distillation* (PSAD), which rebalances supervision across scales and mitigates the gradient imbalance inherent to next-scale decoding. On ImageNet  $256 \times 256$  single-image generation, EVAR achieves substantial reductions in parameters, memory footprint, and end-to-end latency while preserving competitive FID, IS, precision, and recall. In particular, a pruned VAR-d16 model on iPad Pro (M4) attains about  $1.8\times$  speedup in single-image latency with only a modest FID increase. With a practical on-device deployment and evaluation pipeline for single-image, EVAR demonstrates that next-scale VAR can be compressed and executed efficiently on edge hardware, offering a practical step toward visual autoregressive image generation under tight resource constraints.

# EVAR: Edge Visual Autoregressive Models via Principled Pruning

## Supplementary Material

### 537 Appendix Overview

- 538 • **Section A:** Validates generality on LlamaGen-L-256, pre-  
539 serving metrics at 20% sparsity.
- 540 • **Section B:** Reports a 1.8× speedup on iOS devices, ad-  
541 dressing CoreML operator constraints.
- 542 • **Section C:** Compares visual samples, confirming consis-  
543 tency between Linux and mobile deployments.

### 544 A. Generalization on LlamaGen

545 To further investigate the generality of our pruning method,  
546 we conducted experiments on another visual autoregressive  
547 image generation model, LlamaGen. Specifically, we se-  
548 lected LlamaGen-L-256 and applied 20% sparsity pruning,  
549 comparing it against VAR-d16 under the same 20% spar-  
550 sity level. In processing the calibration set for LlamaGen,  
551 since it predicts next tokens, we encoded 256×256 images  
552 to obtain 256 tokens, which served as the pre-encoded cal-  
553 ibration set to guide the structured OBS pruning. Other as-  
554 pects remained consistent with the pruning process for the  
555 VAR model. Since LlamaGen employs next-token genera-  
556 tion rather than next-scale generation, we did not apply  
557 our proposed Progressive Scale-Aware Distillation (PSAD)  
558 method; instead, we used standard fine-tuning. For fairness,  
559 we also applied standard fine-tuning to VAR-d16 as a base-  
560 line comparison.

561 The consolidated results are presented in Table 5. The  
562 findings indicate that our method generalizes effectively to  
563 other autoregressive architectures, reinforcing its broad ap-  
564 plicability and robustness across diverse model families.

Table 5. Performance comparison with and without fine-tuning.

Model	FT	FID ↓	IS ↑	Prec. ↑	Rec. ↑
VAR-d16	w/ FT	3.81	60.43	0.84	0.51
	w/o FT	8.86	48.41	0.74	0.51
LlamaGen-L	w/ FT	3.57	258.3	0.82	0.50
	w/o FT	13.65	43.54	0.71	0.48

### 565 B. Deployment Details and Further Evaluation

566 During the conversion from PyTorch to CoreML, we en-  
567 countered a limitation where the current version of coreml-  
568 tools does not support the bicubic interpolation operator.  
569 Consequently, we replaced it with bilinear sampling for  
570 mobile deployment. This operator switch initially caused  
571 the FID of the VAR-d16 model to increase from 3.55 to  
572 9.98. However, after a brief fine-tuning phase of 20 epochs,

573 the FID recovered to 4.34. While this represents a sys-  
574 tematic increase of approximately 0.78 FID (and a corre-  
575 sponding rise in EVAR from 3.91 to 4.63) compared to  
576 the bicubic baseline, this degradation is strictly attributable  
577 to the toolchain limitation rather than the compression  
578 method. Crucially, the relative quality loss between the  
579 unpruned and pruned models—when both utilize bilinear  
580 sampling—remains within 10%. We anticipate this gap will  
581 vanish as coremltools support evolves.

582 We evaluated inference latency on iPad and iPhone de-  
583 vices using a robust trimmed mean protocol (20 runs after  
584 5 warmups, excluding outliers). Results in Table 6 reveal  
585 two key findings. First, on iPads, our method maintains a  
586 consistent 1.8× speedup. Second, and most notably, both  
587 the full and 20% sparse models crashed on iPhones due to  
588 memory constraints (OOM). Only the 40% sparsity model  
589 operated successfully. This proves that high-ratio compres-  
590 sion is not merely an accelerator but a strict prerequisite for  
591 deploying large generative models on memory-constrained  
592 edge devices.

Table 6. Inference Time Results on Different Devices (ms).

Device	full model	20% Sparsity	40% Sparsity
iPad Pro (M4)	494	439	277
iPad Pro (M2)	791	666	449
iPhone 16 PM	crashed	crashed	580
iPhone 12 Pro	crashed	crashed	1778

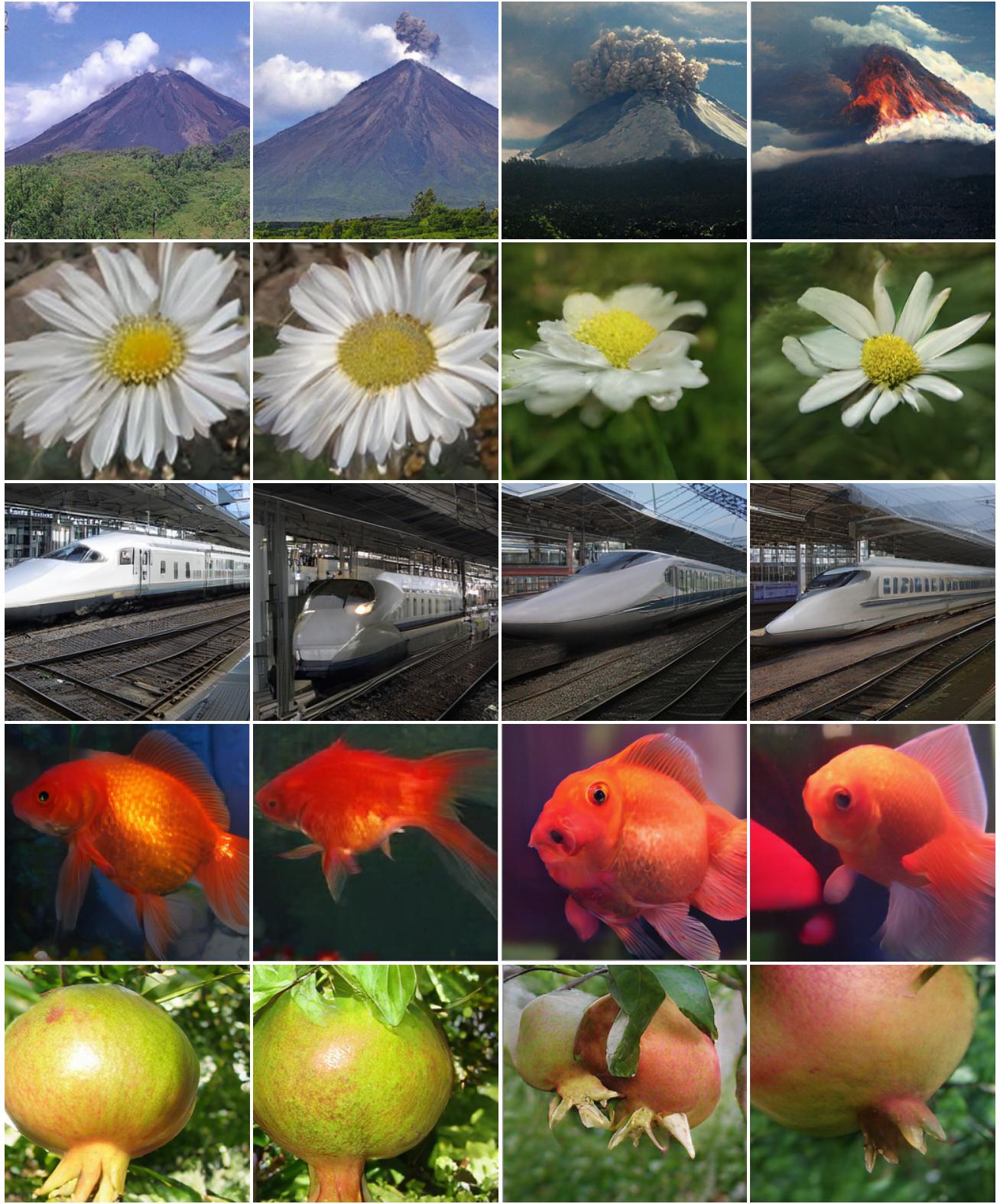
### 593 C. Visual Comparisons

594 Figure 4 presents a visual comparison between the un-  
595 pruned VAR-d16 baseline and our pruned EVAR model  
596 (40% sparsity) under conditional generation settings. To  
597 demonstrate cross-platform consistency, we provide sam-  
598 ples generated in two distinct environments: a standard  
599 Linux workstation using PyTorch and an actual on-device  
600 deployment via CoreML.

601 It is important to note that the images generated on Ap-  
602 ple devices exhibit minor visual deviations compared to the  
603 PyTorch baseline. These differences stem from the model  
604 conversion process and backend constraints. Most notably,  
605 since coremltools does not currently support bicubic sam-  
606 pling, we substituted it with bilinear sampling. Additionally,  
607 there are inherent discrepancies in operator implemen-  
608 tation and fusion strategies between the PyTorch runtime  
609 and the CoreML backend.

610 Despite these constraints, the pruned EVAR model main-  
611 tains high visual fidelity on mobile devices, closely mirror-  
612 ing the original model’s capabilities.

Figure 4. Qualitative comparison of generated samples. We compare the unpruned VAR-d16 baseline executed on a Linux workstation (PyTorch) against our pruned EVAR model with 40% sparsity deployed on mobile devices (CoreML). Despite the aggressive compression and platform constraints, the mobile samples retain high visual fidelity comparable to the baseline.



(a) VAR-d16 on PyTorch

(b) EVAR on PyTorch

(c) VAR-d16 on CoreML

(d) EVAR on CoreML

613 **References**

- [1] Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. Slicept: Compress large language models by deleting rows and columns. In *ICLR*, 2024. 2
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 2
- [3] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. Maskgit: Masked generative image transformer. In *CVPR*, 2022. 2
- [4] Tianyi Chen, Luming Liang, Tianyu Ding, Zhihui Zhu, and Ilya Zharkov. Otov2: Automatic, generic, user-friendly. In *ICLR*, 2023. 2
- [5] Yongwei Chen, Yushi Lan, Shangchen Zhou, Tengfei Wang, and Xingang Pan. Sar3d: Autoregressive 3d object generation and understanding via multi-scale 3d vqvae. In *CVPR*, 2025. 1
- [6] Zigeng Chen, Xinyin Ma, Gongfan Fang, and Xinchao Wang. Collaborative decoding makes visual auto-regressive modeling efficient. In *CVPR*, 2025. 1, 2
- [7] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021. 1, 2, 7
- [8] Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *CVPR*, 2023. 2
- [9] Elias Frantar and Dan Alistarh. Optimal brain compression: A framework for accurate post-training quantization and pruning. In *NeurIPS*, 2022. 2, 4
- [10] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *ICML*, 2023. 2
- [11] Jiatao Gu, Yuyang Wang, Yizhe Zhang, Qihang Zhang, Dinghuai Zhang, Navdeep Jaitly, Joshua M. Susskind, and Shuangfei Zhai. Denoising autoregressive transformers for scalable text-to-image generation. In *ICLR*, 2025. 2
- [12] et al Han, Song. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 2015. 2, 7
- [13] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bit-wise autoregressive modeling for high-resolution image synthesis. In *CVPR*, 2025. 1
- [14] B. Hassibi, D.G. Stork, and G.J. Wolff. Optimal brain surgeon and general network pruning. In *NeurIPS*, 1993. 2
- [15] Haoyu He, Jianfei Cai, Jing Liu, Zizheng Pan, Jing Zhang, Dacheng Tao, and Bohan Zhuang. Pruning self-attentions into convolutional layers in single path. *TPAMI*, 46(5):3910–3922, 2024. 2
- [16] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, 2017. 2
- [17] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, 2017. 2
- [18] Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In *NeurIPS*, 1990. 2
- [19] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image generation using residual quantization. *arXiv preprint arXiv:2203.01941*, 2022. 2, 7
- [20] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *ICLR*, 2017. 2
- [21] Haopeng Li, Jinyue Yang, Kexin Wang, Xuerui Qiu, Yuhong Chou, Xin Li, and Guoqi Li. Scalable autoregressive image generation with mamba. *arXiv preprint arXiv:2408.12245*, 2024. 1, 2, 7
- [22] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024. 2
- [23] Xiang Li, Kai Qiu, Hao Chen, Jason Kuen, Zhe Lin, Rita Singh, and Bhiksha Raj. Controlvar: Exploring controllable visual autoregressive modeling. *arXiv preprint arXiv:2406.09750*, 2024. 1
- [24] Xiang Li, Kai Qiu, Hao Chen, Jason Kuen, Jiuxiang Gu, Bhiksha Raj, and Zhe Lin. Imagefolder: Autoregressive image generation with folded tokens. In *ICLR*, 2025. 1
- [25] Liu Q Ling G, Wang Z. Slimgpt: Layer-wise structured pruning for large language models. In *NeurIPS*, 2024. 2
- [26] Dongyang Liu, Shitian Zhao, Le Zhuo, Weifeng Lin, Yi Xin, Xinyue Li, Qi Qin, Yu Qiao, Hongsheng Li, and Peng Gao. Lumina-mgpt: Illuminate flexible photorealistic text-to-image generation with multimodal generative pretraining. *arXiv preprint arXiv:2408.02657*, 2024. 2
- [27] Enshu Liu, Xuefei Ning, Yu Wang, and Zinan Lin. Distilled decoding 1: One-step sampling of image auto-regressive models with flow matching. In *ICLR*, 2025. 1
- [28] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. In *NeurIPS*, 2024. 2, 7
- [29] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Imagenet: A large-scale hierarchical image database. *CVPR*, 2009. 6
- [30] Tyree S Molchanov P, Mallya A. Importance estimation for neural network pruning. In *CVPR*, 2019. 2, 7
- [31] Aäron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. Conditional image generation with pixelcnn decoders. In *NeurIPS*, 2016. 2
- [32] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2
- [33] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini

- 726 Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob  
727 Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda  
728 Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan  
729 Lowe. Training language models to follow instructions with  
730 human feedback. In *NeurIPS*, 2022. 2
- 731 [34] William Peebles and Saining Xie. Scalable diffusion models  
732 with transformers. *arXiv preprint arXiv:2212.09748*, 2022.  
733 2
- 734 [35] Alec Radford and Karthik Narasimhan. Improving language  
735 understanding by generative pre-training. *Semantic Scholar*,  
736 2018. 2
- 737 [36] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario  
738 Amodei, Ilya Sutskever, et al. Language models are unsu-  
739 pervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 2
- 740 [37] Sucheng Ren, Yaodong Yu, Nataniel Ruiz, Feng Wang,  
741 Alan Yuille, and Cihang Xie. M-var: Decoupled scale-wise  
742 autoregressive modeling for high-quality image generation.  
743 *arXiv preprint arXiv:2411.10433*, 2024. 1
- 744 [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
745 Patrick Esser, and Björn Ommer. High-resolution im-  
746 age synthesis with latent diffusion models. *arXiv preprint*  
747 *arXiv:abs/2112.10752*, 2021. 2
- 748 [39] Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient  
749 second-order approximation for neural network compres-  
750 sion. In *NeurIPS*, 2024. 2
- 751 [40] Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter.  
752 A simple and effective pruning approach for large language  
753 models. In *ICLR*, 2024. 2
- 754 [41] Mingyuan Sun, Zheng Fang, Jiaxu Wang, Junjie Jiang, Delei  
755 Kong, Chenming Hu, Yuetong Fang, and Renjing Xu. Opti-  
756 mal brain apoptosis. In *ICLR*, 2025. 2, 7
- 757 [42] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue  
758 Peng, Ping Luo, and Zehuan Yuan. Autoregressive model  
759 beats diffusion: Llama for scalable image generation. *arXiv*  
760 *preprint arXiv:2406.06525*, 2024. 1, 2, 7
- 761 [43] Haotian Tang, Yecheng Wu, Shang Yang, Enze Xie, Junsong  
762 Chen, Junyu Chen, Zhuoyang Zhang, Han Cai, Yao Lu, and  
763 Song Han. Hart: Efficient visual generation with hybrid au-  
764 toregressive transformer. In *ICLR*, 2025. 1
- 765 [44] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Li-  
766 wei Wang. Visual autoregressive modeling: Scalable image  
767 generation via next-scale prediction. In *NeurIPS*, 2024. 1, 3,  
768 6, 7
- 769 [45] Aäron Van Den Oord, Nal Kalchbrenner, and Koray  
770 Kavukcuoglu. Pixel recurrent neural networks. In *ICML*,  
771 2016. 2
- 772 [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoo-  
773 reit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia  
774 Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2
- 775 [47] Huan Wang, Can Qin, Yulun Zhang, and Yun Fu. Neural  
776 pruning via growing regularization. In *ICLR*, 2021. 2
- 777 [48] Huan Wang, Yulun Zhang, Can Qin, Luc Van Gool, and Yun  
778 Fu. Global aligned structured sparsity learning for efficient  
779 image super-resolution. *TPAMI*, 45(9):10974–10989, 2023.  
780 2
- 781 [49] Jinhong Wang, Jian Liu, Dongqi Tang, Weiqiang Wang,  
782 Wentong Li, Danny Chen, Jintai Chen, and Jian Wu. Scal-
- 783 able autoregressive monocular depth estimation. In *CVPR*,  
784 2025. 1
- 785 [50] Wenhui Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao  
786 Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. PVT  
787 v2: Improved baselines with pyramid vision transformer.  
788 *Comput. Vis. Media*, 8(3):415–424, 2022. 2
- 789 [51] Yuqing Wang, Shuhuai Ren, Zhijie Lin, Yujin Han, Haoyuan  
790 Guo, Zhenheng Yang, Difan Zou, Jashi Feng, and Xihui Liu.  
791 Parallelized autoregressive visual generation. *arXiv preprint*  
792 *arXiv:2412.15119*, 2024. 1, 2, 7
- 793 [52] Ma Xiaoxiao, Zhou Mohan, Liang Tao, Bai Yalong, Zhao  
794 Tiejun, Li Biye, Chen Huaian, and Jin Yi. Star: Scale-wise  
795 text-conditioned autoregressive image generation. *arXiv*  
796 *preprint arXiv:2406.10797*, 2024. 1
- 797 [53] Rui Xie, Tianchen Zhao, Zhihang Yuan, Rui Wan, Wenxi  
798 Gao, Zhenhua Zhu, Xuefei Ning, and Yu Wang. Lite-  
799 var: Compressing visual autoregressive modelling with  
800 efficient attention and quantization. *arXiv preprint*  
801 *arXiv:2411.17178*, 2024.
- 802 [54] Qian Zhang, Xiangzi Dai, Ninghua Yang, Xiang An, Zi-  
803 yong Feng, and Xingyu Ren. Var-clip: Text-to-image gen-  
804 erator with visual auto-regressive modeling. *arXiv preprint*  
805 *arXiv:2408.01181*, 2024. 1
- 806 [55] Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala,  
807 Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe  
808 Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Pre-  
809 dict the next token and diffuse images with one multi-modal  
810 model. *arXiv preprint arXiv:2408.11039*, 2024. 2