

SparVAR: Exploring Sparsity in Visual AutoRegressive Modeling for Training-Free Acceleration

Zekun Li^{1,2,3} Ning Wang^{1,2,3} Tongxin Bai^{3*} Changwang Mei^{1,4}
 Peisong Wang^{1,2*} Shuang Qiu⁵ Jian Cheng^{1,2*}

¹Institute of Automation, Chinese Academy of Sciences,

²School of Artificial Intelligence, University of Chinese Academy of Sciences,

³Beijing Academy of Artificial Intelligence,

⁴Nanjing University of Science and Technology,

⁵City University of Hong Kong.



Figure 1. Our **SparVAR** achieves efficient acceleration while preserving high-frequency details consistent with the baseline [12], whereas prior methods introduce visible artifacts and texture loss. The bottom metrics denotes GenEval / PSNR. Zoom in for fine-detail comparison.

Abstract

Visual AutoRegressive (VAR) modeling has garnered significant attention for its innovative next-scale prediction paradigm. However, mainstream VAR paradigms attend to all tokens across historical scales at each autoregressive step. As the next scale resolution grows, the computational complexity of attention increases quartically with

resolution, causing substantial latency. Prior accelerations often skip high-resolution scales, which speeds up inference but discards high-frequency details and harms image quality. To address these problems, we present **SparVAR**, a training-free acceleration framework that exploits three properties of VAR attention: (i) **strong attention sinks**, (ii) **cross-scale activation similarity**, and (iii) **pronounced locality**. Specifically, we dynamically predict the sparse attention pattern of later high-resolution scales from a sparse

*Corresponding authors.

decision scale, and construct scale self-similar sparse attention via an efficient index-mapping mechanism, enabling high-efficiency sparse attention computation at large scales. Furthermore, we propose cross-scale local sparse attention and implement an efficient block-wise sparse kernel, which achieves $> 5\times$ faster forward speed than FlashAttention. Extensive experiments demonstrate that the proposed SparseVAR can reduce the generation time of an 8B model producing 1024×1024 high-resolution images to the **1s, without skipping the last scales**. Compared with the VAR baseline accelerated by FlashAttention, our method achieves a $1.57\times$ speed-up while preserving almost all high-frequency details. When combined with existing scale-skipping strategies, SparseVAR attains up to a $2.28\times$ acceleration, while maintaining competitive visual generation quality. Code is available at <https://github.com/CAS-CLab/SparVAR>.

1. Introduction

Recent advances in visual generative modeling have been dominated by two major paradigms: diffusion models [2, 3, 15, 20, 31, 32, 36, 61] and autoregressive (AR) architectures [8, 24, 26, 40, 50, 56]. Diffusion models achieve remarkable image quality via iterative denoising but suffer from slow generation due to long sampling chains, whereas AR models generate visual tokens sequentially [57], typically offering higher sampling efficiency [30, 51]. Among them, visual autoregressive (VAR) modeling [12, 41, 42] has recently emerged as a promising alternative by introducing a novel next-scale prediction paradigm: instead of generating individual token at each step, VAR predicts all tokens of the next scale in parallel, progressively refining higher-resolution residuals [13] and enabling more efficient and scalable inference.

Despite its strong potential, current VAR frameworks still suffer from high inference latency when generating high-resolution images. In next-scale prediction, tokens at the current scale must attend to all tokens from previous scales to maintain structural consistency, causing the attention complexity [47] to grow roughly quartically ($\mathcal{O}(n^2) \rightarrow \mathcal{O}(n^4)$) with image resolution and making the last two large-scale steps account for about 60% of the total runtime [11]. Meanwhile, the accumulation of key-value (KV) caches across historical scales drastically increases GPU memory usage—an 8B VAR model [12] requires nearly 60 GB to generate 1024×1024 images [22]—thereby limiting batch inference, overall throughput, and deployment scalability.

To alleviate this bottleneck, recent acceleration techniques for VAR focus primarily on **skipping late high-resolution scales**. By omitting the last two or three scales, methods such as FastVAR [11] and SkipVAR [21] substan-

tially reduce computation and inference latency while maintaining semantic correctness. However, this improvement comes **at the cost of losing fine-grained high-frequency details**, as the last scales are responsible for texture refinement, spatial sharpness and even detailed content generation. As shown in Fig. 1, although high-level metrics such as GenEval [10] score may still indicate semantic alignment, low-level metrics—including Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) [52], and Learned Perceptual Image Patch Similarity (LPIPS) [60]—clearly reveal a substantial drop in visual fidelity compared with the baseline model [12]. Even adaptive strategies like SkipVAR, which determine skipping dynamically per-sample, still yield noticeably lower low-level scores. We further observe that in multi-object or fine-texture scenarios, skipping high-resolution scales often leads to **missing texture patterns** and **structural distortions** in the generated images.

These limitations raise a question: *Can we achieve effective acceleration for VAR without skipping any scales?*

To answer this question, we conduct a systematic analysis of the attention activation patterns in Infinity [12], a representative text-to-image VAR model, and reveal three consistent properties across layers, heads, and scales: **1) Strong Attention Sinks**. As shown in Fig. 2(a), we observe that a small subset of early-scale tokens consistently attracts high attention weights, functioning as “global anchors”. We further investigate the influence of historical scale KV caches on the final generation results: as illustrated in Fig. 3, when retaining only the first 4 \sim 5 scales of the KV cache, the model is still able to reconstruct accurate object layouts and global structures. **2) Cross-Scale Activation Similarity**. As shown in Fig. 2(b), for instance, the attention of tokens within scale 10 (highlighted by the red box) closely resembles that of scales 11 and 12. Moreover, in the cross-scale regions where tokens from the current scale interact with those from historical scales (yellow box), similar attention structures can also be observed. This similarity indicates that the attention of the last scales can be effectively predicted from earlier scales, providing a solid foundation for efficient sparse computation at high-resolutions scales. **3) Pronounced Spatial Locality**. As the resolution of next-scale increases, the attention patterns become increasingly concentrated within local spatial neighborhoods of the same or adjacent scales. In particular, this locality manifests as cross-scale diagonal activation patterns in the attention maps. These systematic findings reveal a high degree of redundancy in the attention activation patterns of VAR, both intra- and cross-scale, inspiring us to develop a training-free sparse acceleration framework for efficient inference.

Building upon these observations, we introduce **SparVAR**, a **training-free acceleration framework** that exploits the sparsity of VAR attention. Specifically, we

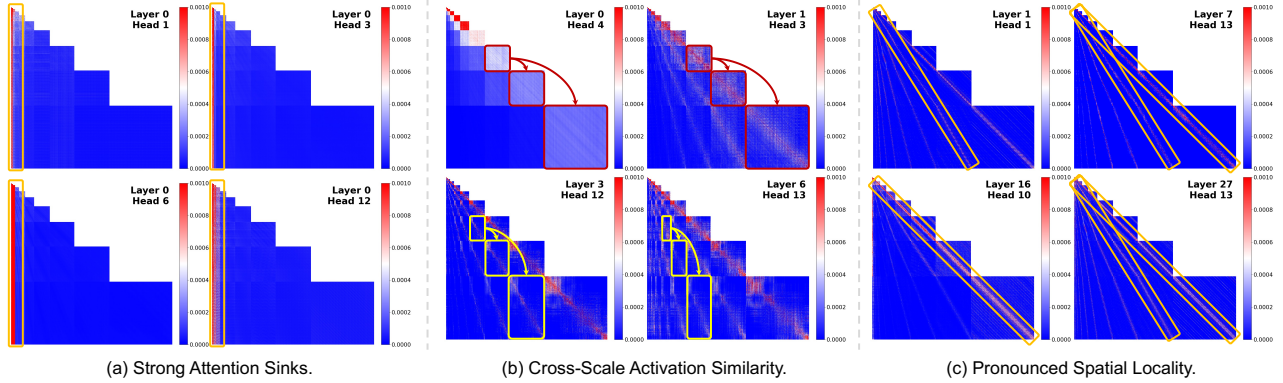


Figure 2. Visualization of attention activation patterns in the Infinity [12] across different layers and heads. (a) Strong Attention Sinks: early-scale tokens consistently attract large attention weights, serving as global anchors that dominate image structure formation. (b) Cross-Scale Activation Similarity: corresponding sub-blocks across adjacent scales exhibit similar activation distributions, indicating redundant attention patterns that can be transferred across scales. (c) Pronounced Spatial Locality: at higher scales, attention becomes increasingly concentrated along local spatial bands, revealing strong locality both within and between neighboring scales.



Figure 3. The manifestation of attention sinks in the KV cache.

introduce two plug-and-play efficient attention modules. 1) **Cross-Scale Self-Similar Sparse Attention (CS^4A)**, which dynamically predicts the sparse attention patterns of high-resolution scales from a sparse decision scale via an efficient cross-scale sparse index mapping. SparVAR transfers sparsity across scales and performs attention computation only over the selected sparse KV Cache at large scales, significantly improving efficiency. 2) **Cross-scale local sparse attention ($CSLA$)**, which reinforces spatial locality via an optimized **Block-wise Local Sparse Kernel** ($> 5\times$ faster than Flash-attention [5, 6] in forward), achieving substantial computational savings while preserving fine-grained details. Together, these components enable SparVAR to accelerate inference without skipping any scales, allowing an 8B model to generate 1024×1024 images at second-level latency with a $1.57\times$ speed-up. Most importantly, SparVAR maintains semantic consistency and achieves nearly identical high-fidelity image quality compared with the baseline model.

Our main contributions are summarized as follows:

- We systematically analyze the attention activation patterns in pretrained VAR models and observe three key phenomena—strong attention sinks, cross-scale activation similarity, and pronounced spatial locality—which directly motivate the exploration of attention sparsity in VAR.
- Based on these findings, we design two plug-and-play sparse attention modules and implement customized block-wise sparse kernels that achieve significant accel-

eration while preserving fine-grained generation details.

- The proposed SparVAR is a training-free inference acceleration framework that exploits cross-scale attention sparsity in VAR, achieving up to $1.57\times$ faster inference without sacrificing high-frequency fidelity.

2. Related Work

Autoregressive Visual Generation. Early autoregressive (AR) image generators flatten a 2D image into a 1D raster-scan sequence and predict the next pixel [37, 44, 45], which becomes prohibitively expensive at high resolutions. To improve scalability, subsequent work adopts vector quantization [35, 46] to convert image patches into discrete tokens and pairs them with Transformers [47] in a next-token paradigm [8, 26, 40]. Leveraging large language model scaling laws [14, 19, 34], token-based AR achieves image quality competitive with state-of-the-art diffusion models [1–3, 9, 32, 36]. However, strictly sequential token generation still incurs substantial latency and raster flattening weakens 2D inductive bias, making efficient high-resolution synthesis challenging. Recently, visual autoregressive (VAR) modeling [12, 41, 42] adopts a next-scale paradigm: at each AR step the model predicts all tokens of the next resolution scale in parallel, generating images in a coarse-to-fine manner. Yet at each step VAR attention must attend to all tokens from previous scales, so the computational and memory costs grow rapidly with resolution. This scalability issue severely limits the throughput and responsiveness of large VAR models for high-resolution generation (e.g., 1024×1024 and beyond).

Training-Free Acceleration for Visual Generation. Training-free acceleration improves inference efficiency of pretrained models without additional fine-tuning by exploiting their intrinsic structure. For diffusion models,

such methods are well developed: improved samplers and scheduling [17, 23, 25, 28] reduce denoising iterations or adapt the trajectory, while feature reuse and caching [4, 18, 27, 29, 53, 59, 62] leverage slowly-varying activations across timesteps. Chipmunk [39] further exploits activation sparsity by recomputing only the fastest-changing activations, showing that *attention activation pattern reuse is a powerful training-free prior*. In contrast, training-free acceleration for VAR remains underexplored. Fast-VAR [11] dynamically prunes tokens and skips the last high-resolution scales to obtain large speedups, but causes noticeable loss of high-frequency details and structural artifacts. ScaleKV [22] and HACK [33] compress the KV cache and improve memory efficiency but yield limited runtime gains due to expensive key-value selection, while SkipVAR [21] uses a lightweight discriminator for sample-adaptive scale skipping at the cost of retraining and quality degradation. Our SparVAR is a training-free acceleration framework that leverages VAR-specific attention activation sparsity for efficient computation, achieving over $1.5\times$ speedup without skipping the last scales.

3. Method

3.1. Background

Visual AutoRegressive (VAR) modeling [42] redefines the traditional autoregressive paradigm by shifting from a *next-token* prediction to a *next-scale* prediction paradigm. Instead of generating one token sequentially, each autoregressive step in VAR produces an entire token map at a specific resolution scale, progressively refining the image from coarse to fine. Given an image feature map $\mathbf{F} \in \mathbb{R}^{H \times W \times D}$, VAR quantizes it into K multi-scale residual token maps

$$\mathcal{R} = \{\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_K\},$$

where \mathbf{R}_k denotes the token map at scale k . The resolution of \mathbf{R}_k is $h_k \times w_k$ and it grows larger gradually from $k = 1 \rightarrow K$. The size of last scale token map is $(h_K \times w_K) = (H \times W)$. Based on the text prompt t embedding $\Psi(t)$ and \mathcal{R} , the joint distribution is factorized autoregressively as:

$$p(\mathcal{R}) = \prod_{k=1}^K p(\mathbf{R}_k | \mathbf{R}_1, \dots, \mathbf{R}_{k-1}, \Psi(t)),$$

where $\Psi(t)$ typically is projected to the first scale with size 1×1 . During inference, the model predicts the next-scale residual \mathbf{R}_k in parallel conditioned on all previous scales.

Let $N_k = h_k w_k$ denote the number of tokens at scale k , then the total number of tokens across all preceding scales including scale k is given by $N_{\leq k} = \sum_{i=1}^k N_i$. Each Transformer block employs causal cross-scale attention, allowing tokens of scale k to attend to tokens from all previous scales

$\mathbf{R}_{\leq k}$:

$$\mathbf{O}_k = \text{Softmax} \left(\frac{\mathbf{Q}_k \mathbf{K}_{\leq k}^\top}{\sqrt{d}} \right) \mathbf{V}_{\leq k}, \quad (1)$$

where $\mathbf{Q}_k \in \mathbb{R}^{B \times H \times N_k \times D}$ denotes the query projection from the current scale tokens, and $\mathbf{K}_{\leq k}, \mathbf{V}_{\leq k} \in \mathbb{R}^{B \times H \times N_{\leq k} \times D}$ represent the concatenation of key and value projections from both the current and previous scales along the sequence dimension. This enables coarse-to-fine dependency modeling, but as the scale spatial resolution increases, the number of tokens grows quadratically, resulting in $\mathcal{O}(n^4)$ attention complexity and substantial memory overhead. Therefore, designing efficient cross-scale attention is crucial for scaling VAR inference to high resolutions—this motivates our SparVAR.

3.2. Cross-Scale Self-Similar Sparse Attention

Motivation. To systematically analyze the pretrained text-to-image VAR model [12], we visualize the attention maps at each scale and vertically concatenate them in ascending order of spatial resolution. The resulting visualization, shown in Fig. 2, clearly reveals how attention activation patterns evolve across scales. From Fig. 2(b), it is evident that the attention activation patterns across different scales exhibit pronounced cross-scale similarity—the activation at scale k closely mirrors that of the preceding scale $k-1$, particularly in the corresponding subregions when the attention map is partitioned by source scales. Formally, for the full attention map $\mathbf{A}^{(k)} \in \mathbb{R}^{B \times H \times N_k \times N_{\leq k}}$ at scale k , we divide it into k column-blocks according to the source scales:

$$\mathbf{A}^{(k)} = [\mathbf{A}^{(k,1)} \quad \mathbf{A}^{(k,2)} \quad \dots \quad \mathbf{A}^{(k,k)}], \quad (2)$$

where each block $\mathbf{A}^{(k,i)} \in \mathbb{R}^{N_k \times N_i}$ corresponds to the attention from the current scale k to the i -th previous scale ($i \leq k$). We empirically observe a strong similarity between consecutive scales:

$$\mathbf{A}^{(k,i)} \approx \text{Upsample}(\mathbf{A}^{(k-1,i-1)}), \quad \forall i \geq 2,$$

suggesting that the activation pattern of the current scale can be predicted by upsampling that of the preceding scale. This motivates us to reuse the sparse activation pattern identified at a mid-scale to predict sparse attention for larger scales, thereby avoiding redundant dense computations.

Sparse Decision Scale. We define a **sparse decision scale** S as the optimal mid-scale whose dense attention provides a representative sparse pattern for predicting the sparse activations of later scales. Mathematically, S is selected to minimize the discrepancy between its predicted sparse pattern (after cross-scale mapping) and the dense attention of the target high-resolution scale K :

$$S = \arg \min_{S \in 1, \dots, K-1} \mathcal{L}_{\text{perf}}(\mathcal{M}_{S \rightarrow K}(\Omega^{(S)}), \Omega_{\text{dense}}^{(K)}),$$

where $\Omega^{(s)} = \text{TopK}(\mathbf{A}^{(s)})$ denotes the sparse activation set extracted from the dense attention of scale S , $\mathcal{M}_{S \rightarrow K}(\cdot)$ represents the cross-scale sparse index mapping function, and $\mathcal{L}_{\text{perf}}$ measures the generation fidelity between the predicted and dense patterns. In practice, we fix S to a mid-scale (e.g., 10 scale), which—as shown in our ablation studies—achieves near-baseline generation performance while substantially reducing computational cost. The sparse indices $\text{inds}^{(S)}$ extracted at this scale serve as the foundation for all subsequent sparse inference.

Specifically, at scale S , we compute full scaled dot-product attention following Eq. 1:

$$\mathbf{P}^{(S)} = \text{Softmax} \left(\frac{\mathbf{Q}^{(S)} \mathbf{K}^{(S)\top}}{\sqrt{d}} \right), \mathbf{O}_{\text{dense}}^{(S)} = \mathbf{P}^{(S)} \mathbf{V}^{(S)}.$$

We partition $\mathbf{P}^{(S)}$ along the query axis into $G_S = N_S/C$ contiguous query blocks of size C

$$\mathbf{P}^{(S)} \in \mathbb{R}^{B \times H \times N_S \times N_{\leq S}} \rightarrow \text{reshape}(B, H, \frac{N_S}{C}, C, N_{\leq S}),$$

and sum over each block to obtain the column-sum tensor [39]:

$$\mathbf{D}_{b,h,i,j}^{(S)} = \sum_{q=ic}^{(i+1)c-1} \mathbf{P}_{b,h,q,j}^{(S)},$$

where $\mathbf{D}^{(S)} \in \mathbb{R}^{B \times H \times G_S \times N_{\leq S}}$ encodes the total attention weight each block assigns to key j . A Top- k operation then identifies the most attended key indices:

$$\text{inds}_{b,h,i,:k}^{(S)} = \text{TopK}_k(\mathbf{D}_{b,h,i,:}^{(S)}).$$

Finally, we define the dense attention cache by subtracting the sparse component:

$$\mathbf{O}_{\text{cache}}^{(S)} = \mathbf{O}_{\text{dense}}^{(S)} - \text{Softmax} \left(\mathbf{Q}^{(S)} \mathbf{K}_{\text{inds}}^{(S)\top} \right) \mathbf{V}_{\text{inds}}^{(S)},$$

which captures the residual part of the dense output after removing the top sparse activations.

Cross-Scale Sparse Computation. For subsequent high-resolution scales $k > S$, we reuse the sparse indices $\text{inds}^{(S)}$ by mapping them across scales to obtain the target sparse indices $\text{inds}^{(k)}$. We design an efficient **cross-scale sparse index mapping**, which projects $\text{inds}^{(S)}$ onto the spatial grid of scale k . We compute the mapped indices as:

$$\text{inds}^{(k)} = \left\lfloor \frac{N_k}{N_S} \times \text{inds}^{(S)} \right\rfloor,$$

optionally adding a sink region \mathcal{S} to guarantee coverage of global anchors:

$$\text{inds}^{(k)} \leftarrow \text{Concat}(\mathcal{S}, \text{inds}^{(k)}).$$

This mapping preserves the hierarchical correspondence between scales, ensuring that attention sparsity remains coherent across resolutions. Further details can be found in the appendix. Then the sparse update at scale k is then efficiently computed as:

$$\Delta \mathbf{O}^{(k)} = \text{Softmax} \left(\mathbf{Q}^{(k)} \mathbf{K}_{\text{inds}^{(k)}}^{(k)\top} \right) \mathbf{V}_{\text{inds}^{(k)}}^{(k)},$$

$$\mathbf{O}^{(k)} \approx \mathbf{O}_{\text{cache}}^{(k)} + \Delta \mathbf{O}^{(k)},$$

where $\mathbf{O}_{\text{cache}}^{(k)}$ denotes the upsampled cache prediction from scale S . As shown in Tab. 5, the $\mathbf{O}_{\text{cache}}^{(k)}$ term further enhances the high-frequency visual fidelity of the accelerated model with minimal latency.

3.3. Cross-Scale Local Sparse Attention

Motivation. Building upon the cross-scale sparse pattern reuse introduced in CS^4A , we further observe that in deep layers of pretrained VAR models, attention activation at later high-resolution scales primarily concentrates on two regions: (i) **attention sink** formed by early scales, and (ii) **spatial locality** both within and across adjacent scales. Hence, instead of relying on mapped sparse indices from decision scale, directly exploiting these structural sparsity priors allows for greater acceleration benefits. Moreover, recent advances [5–7, 38, 43, 49] in AI system for efficient operator optimization demonstrate that defining sparsity at the **block granularity** fully leverages modern GPU hardware features, achieving extreme operator acceleration. Inspired by these insights, we propose **Cross-Scale Local Sparse Attention (CSLA)**, which constructs an efficient **block-wise sparse attention kernel** guided by sink and locality priors, further accelerating VAR inference while maintaining high-fidelity generation.

Cross-Scale Local Mapping. To exploit the observed spatial locality in VAR attention activation, we define a cross-scale local mapping function that restricts each query token to attend only to its nearby key-value tokens across scales. For each query token \mathbf{q} at scale k , we first identify its corresponding spatial location $(x_{\mathbf{q}}^{(k)}, y_{\mathbf{q}}^{(k)})$ and compute its aligned coordinates in the key grid of a historical scale $h \in \{0, 1, \dots, k-1, k\}$:

$$(\tilde{x}_{\mathbf{q}}^{(k \rightarrow h)}, \tilde{y}_{\mathbf{q}}^{(k \rightarrow h)}) = \left(\text{round} \left(\frac{x_{\mathbf{q}}^{(k)}}{H_k} H_h \right), \text{round} \left(\frac{y_{\mathbf{q}}^{(k)}}{W_k} W_h \right) \right),$$

where (H_h, W_h) denotes the spatial resolution of historical scale h . Each query thus attends to a local neighborhood centered at $(\tilde{x}_{\mathbf{q}}^{(k \rightarrow h)}, \tilde{y}_{\mathbf{q}}^{(k \rightarrow h)})$ with a per-scale window radius r_h , forming a token-wise local sparse mask

$$\mathbf{M}_{\text{local}}^{(k,h)}(\mathbf{q}, \mathbf{k}) = 1 \left(|\tilde{x}_{\mathbf{q}}^{(k \rightarrow h)} - x_{\mathbf{k}}^{(h)}| \leq r_h \wedge |\tilde{y}_{\mathbf{q}}^{(k \rightarrow h)} - y_{\mathbf{k}}^{(h)}| \leq r_h \right).$$

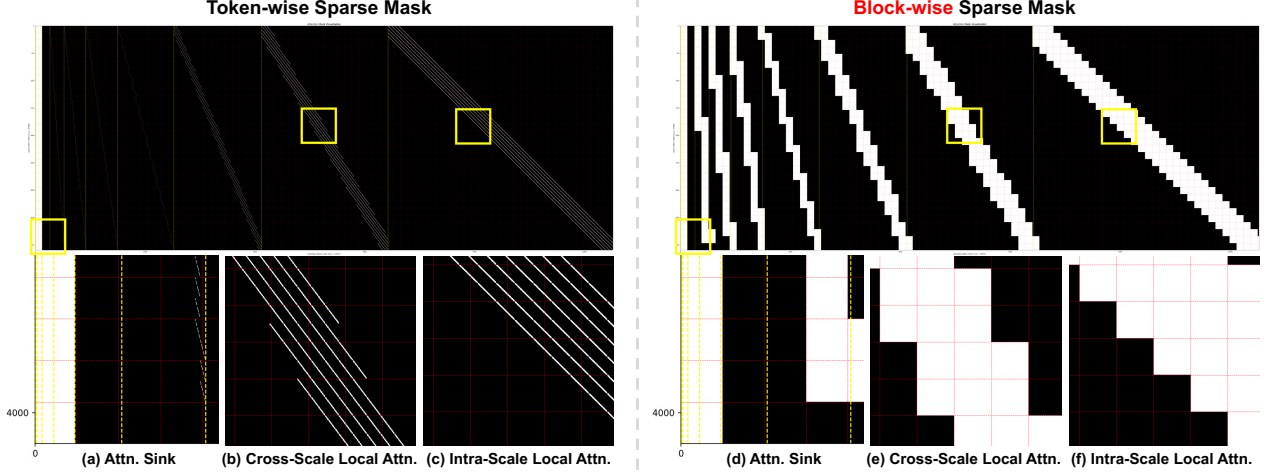


Figure 4. Visualization of the Cross-Scale Local Sparse Attention (CSLA) masks. This example shows the last scale attention map in Infinity ($q_{len} = 4096, kv_{len} = 10521$). The left shows the token-wise sparse mask, and the right shows the corresponding block-wise version after applying the block aggregation in Eq. 3. Red dashed grids denote 128×128 blocks, while yellow dashed lines mark the attention partition boundaries described in Eq. 2. Zoom in for fine-detail visualization.

Intuitively, $M_{\text{local}}^{(k,h)}(\mathbf{q}, \mathbf{k})$ identifies whether key \mathbf{k} lies within the local neighborhood of query \mathbf{q} after rescaling between the two scales. Meanwhile, a small number of early-scale tokens, denoted by the “sink region” \mathcal{S} , remain fully visible to all queries to preserve global context:

$$M_{\text{sink}}^{(k)}(\mathbf{q}, \mathbf{k}) = 1[\mathbf{k} \in \mathcal{S}].$$

Finally, we integrate the above local and sink masks into a unified token-wise sparse mask

$$M^{(k)} = M_{\text{sink}}^{(k)} \vee \bigvee_{i=1}^k M_{\text{local}}^{(k,h)},$$

where \vee denotes the logical OR operation across multiple historical scales.

Block-wise Sparse Mask. To align with GPU-friendly block-sparse computations, we partition the query axis and the key axis uniformly into blocks of size B , where $u \in \{0, \dots, \lceil N_k/B \rceil - 1\}$ denotes the block index along the query dimension and $v \in \{0, \dots, \lceil N_{\leq k}/B \rceil - 1\}$ represents the block index along the key dimension. Furthermore, we construct an efficient block-wise sparse mask:

$$B^{(k)}(u, v) = \bigvee_{\mathbf{q} \in B_u} \bigvee_{\mathbf{k} \in B_v} M^{(k)}(\mathbf{q}, \mathbf{k}), \quad (3)$$

where each block groups $B \times B$ tokens.

Fig. 4 provides an example of the sparse mask employed for attention computation at the last scale. This conversion ensures that every active block represents a dense local window and inactive ones are completely pruned, producing a structured sparsity pattern well aligned with FlashAttention-style tiling. We implement CSLA based

Table 1. Forward speed of sparse attention kernels in a setup aligned with Infinity-8B last scale inference configuration (bf16, $q_{len} = 4096, kv_{len} = 10521, d_{head} = 128, heads = 24$). Test on H100 GPU with CUDA 12.8 and PyTorch 2.7.1.

Operation	Implementation	Sparsity (%)	Latency (ms) ↓	Speedup ↑
FlashAttention2	CUDA	0.00%	3.0231	1.00×
F.sdp + sparse mask	PyTorch	87.86%	8.2276	0.37×
CSLA (w/ block size 64)	FlexAttention	83.50%	0.8396	3.60×
CSLA (w/ block size 128)	FlexAttention	83.46%	0.5392	5.61×

on top of the FlexAttention [7] framework, with quantitative inference results summarized in Tab. 1. Under a block size of 128, our CSLA kernel achieves a **5.61×** forward speed than FlashAttention on the last scale, and a remarkable **15.26×** acceleration compared to the naïve token-wise sparse attention baseline.

4. Experiments

4.1. Experimental Setup

Base Models. We apply our proposed **SparVAR** to the representative VAR-based text-to-image baselines [12], Infinity-2B and -8B, to verify the generality of our approach across different model parameter scales. For a fair comparison, all hyperparameters and model configurations are kept identical to the official Infinity implementations, except for the introduced sparse attention modules. All our speed measurements and evaluation experiments are conducted on NVIDIA H100 GPUs. For more implementation details, please refer to the appendix.

Evaluation Metrics. We evaluate the proposed method from two perspectives: generation quality and inference ef-

Table 2. Quantitative comparison on efficiency and quality on 1024×1024 GenEval [10] benchmarks. Note that the efficiency of Infinity baselines are tested under FlashAttention.

Methods	Acceleration				GenEval Score				Quality Metrics		
	#Scales ↓	Speedup ↑	Latency ↓	Throughput ↑	Two Obj. ↑	Position ↑	Color Attri. ↑	Overall ↑	PSNR ↑	SSIM ↑	LPIPS ↓
<i>w/o skip scales</i>											
Infinity-2B [12]	13	1.00×	0.96s	1.04it/s	0.864	0.450	0.555	0.736	-	-	-
FastVAR [11]	13	1.01×	0.95s	1.05it/s	0.831	0.438	0.520	0.712	15.618	0.632	0.428
ScaleKV [22]	13	0.70×	1.37s	0.73it/s	0.854	0.453	0.548	0.732	23.370	0.820	0.192
SparVAR (Ours)	13	1.38×	0.69s	1.44it/s	0.856	0.458	0.563	0.738	28.882	0.914	0.097
<i>w/ skip last 2 scales</i>											
FastVAR [11]	11	1.33×	0.72s	1.39it/s	0.831	0.405	0.575	0.716	12.867	0.544	0.538
SparVAR (Ours)	11	1.70×	0.56s	1.77it/s	0.846	0.413	0.588	0.725	15.272	0.596	0.474
<i>w/ dynamic skip last scales</i>											
SkipVAR-2B [21]	10 ~ 13	1.33×	0.72s	1.39it/s	0.854	0.418	0.565	0.730	17.650	0.678	0.391
SparVAR (Ours)	10 ~ 13	1.58×	0.61s	1.65it/s	0.854	0.420	0.568	0.732	17.689	0.679	0.389
<i>w/o skip scales</i>											
Infinity-8B [12]	13	1.00×	1.65s	0.61it/s	0.899	0.610	0.695	0.798	-	-	-
FastVAR [11]	13	1.14×	1.45s	0.69it/s	0.917	0.610	0.680	0.792	17.403	0.630	0.333
ScaleKV [22]	13	0.67×	2.45s	0.41it/s	0.886	0.618	0.690	0.793	23.423	0.803	0.153
SparVAR (Ours)	13	1.57×	1.05s	0.95it/s	0.897	0.631	0.683	0.796	29.481	0.920	0.073
<i>w/ skip last 2 scales</i>											
FastVAR [11]	11	1.79×	0.92s	1.09it/s	0.886	0.600	0.685	0.790	15.265	0.533	0.421
SparVAR (Ours)	11	2.28×	0.72s	1.38it/s	0.891	0.608	0.700	0.800	17.096	0.579	0.374
<i>w/ dynamic skip last scales</i>											
SkipVAR-8B [21]	10 ~ 13	1.71×	0.97s	1.04it/s	0.884	0.613	0.675	0.789	17.980	0.641	0.337
SparVAR (Ours)	10 ~ 13	2.05×	0.80s	1.24it/s	0.894	0.625	0.695	0.796	19.614	0.694	0.296

iciency. For generation quality, we consider both *high-level* and *low-level* metrics. High-level evaluation focuses on semantic alignment and human preference, measured using four popular benchmarks—GenEval [10], DPG-Bench [16] (in appendix), ImageReward [58] and HPSv2.1 [54]. For low-level evaluation, we take the baseline model’s outputs as the reference and employ PSNR, SSIM, and LPIPS to quantitatively assess the preservation of high-frequency textures and visual fidelity. It is important to note that directly skipping late high-resolution scales can bring substantial acceleration to VAR models. Therefore, we evaluate generation quality and inference efficiency under three settings: (i) *without scale skipping*, (ii) *skipping the last two scales* [11], and (iii) *sample dynamic skipping* [21]. This comprehensive and fair evaluation setup enables a more thorough comparison across different acceleration strategies.

4.2. Main Results

Comparison on GenEval. We systematically evaluate SparVAR on 1024×1024 text-to-image generation, comparing against state-of-the-art VAR acceleration methods across acceleration, GenEval scores, and low-level fidelity metrics. The specific results are presented in Tab. 2. Without skipping any scales, SparVAR achieves **1.38×** and **1.57×** speedups on Infinity-2B and -8B respectively, while maintaining high-level semantic alignment nearly identical to the baseline. Notably, it reduces the generation time

of the **8B model** to the **1-second** regime for producing a **1024×1024 image**. Compared with FastVAR [11] and ScaleKV [22], SparVAR not only accelerates inference but also preserves high-frequency details far more effectively, achieving substantially better low-level metrics. Specifically, SparVAR attains **PSNR approaching 30**, **SSIM above 0.9**, and **LPIPS below 0.1**, indicating that the proposed cross-scale sparse attention successfully compresses redundant computation while retaining the critical token interactions necessary for fine-grained structure generation. When combined with skip-scale strategies, SparVAR further boosts acceleration to **1.70×** and **2.28×**, increasing the inference throughput of the **2B model** to **1.77 it/s**. Remarkably, even under the extreme acceleration scenario with scale skipping, SparVAR still preserves high-frequency textures effectively, thus maintaining its leading performance on low-level metrics.

Human Preference Evaluation. We further evaluate the generation quality of SparVAR and other methods using human preference benchmarks [54, 58]. As shown in Tab. 3, both 2B and 8B model scales, SparVAR without scale skipping achieves comparable human preference scores to the baseline, even attaining a score of **31.01** at 8B—**surpassing the baseline**. This confirms our cross-scale sparse attention preserves fine semantic alignment and visual realism. When combined with both scale-skipping setup, our SparVAR still demonstrates remarkable robustness under aggressive acceleration. We attribute this advantage to the suppression

of attention noise through sparsity and the enhanced spatial locality that contributes to more consistent generation. For additional qualitative results, please refer to the appendix.

Table 3. Human preference evaluation on HPSv2.1 [54] and ImageReward [58].

Methods	HPSv2.1			ImageReward
	Photo \uparrow	Concept-art \uparrow	Average \uparrow	Avg. Score \uparrow
<i>w/o skip scales</i>				
Infinity-2B [12]	29.45	30.47	30.55	0.9443
FastVAR [11]	28.59	29.81	29.85	0.9327
ScaleKV [22]	29.35	30.31	30.41	0.9327
SparVAR (Ours)	29.41	30.44	30.53	0.9416
<i>w/ skip last 2 scales</i>				
FastVAR [11]	28.80	29.91	29.97	0.9191
SparVAR (Ours)	29.33	30.31	30.41	0.9201
<i>w/ dynamic skip last scales</i>				
SkipVAR-2B [21]	29.31	30.38	30.47	0.9399
SparVAR (Ours)	29.33	30.39	30.47	0.9455
<i>w/o skip scales</i>				
Infinity-8B [12]	29.49	31.28	31.00	1.0529
FastVAR [11]	29.13	30.85	30.62	1.0380
ScaleKV [22]	29.37	31.10	30.82	1.0350
SparVAR (Ours)	29.47	31.34	31.01	1.0533
<i>w/ skip last 2 scales</i>				
FastVAR [11]	28.82	30.38	30.23	1.0276
SparVAR (Ours)	29.19	30.84	30.59	1.0377
<i>w/ dynamic skip last scales</i>				
SkipVAR-8B [21]	29.09	30.90	30.64	1.0297
SparVAR (Ours)	29.15	30.97	30.69	1.0305

4.3. Ablation Studies

Sparse Decision Scale. To determine the optimal sparse decision scale, we conduct both qualitative and quantitative ablation studies. We first analyze how selecting different scales as the sparse decision scale affects the final generation quality. As shown in Fig. 5, using earlier scales (e.g., 6 or 7) introduces severe artifacts, texture corruption, and global blurring, whereas mid-to-late scales consistently preserve high-frequency details. This behavior is tightly linked to the spatial resolution of VAR scales: earlier scales have limited resolution and exhibit coarse, highly dispersed attention patterns that share low similarity with the sparse patterns of later high-resolution scales. Consequently, the cross-scale mapping becomes less accurate, leading to degraded generation quality. We further perform systematic quantitative experiments to identify the most suitable sparse decision scale and the optimal Top-K range. As shown in Fig. 6, scale 10 emerges as the best trade-off between quality and inference speed. Across all settings, a Top-K around 0.2 provides a stable and effective balance between sparsity and fidelity.

Local Sparsity. Tab. 4 presents an ablation on the attention sink range and the window size configuration in CSLA. We observe that retaining early-scale tokens (i.e. at least

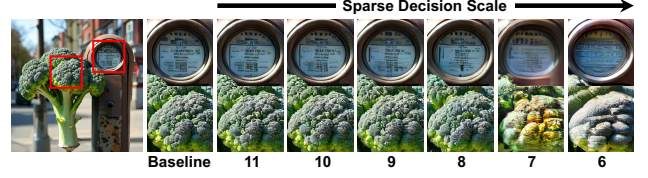


Figure 5. The impact of selecting sparse decision scales on generation results. Zoom in for fine-detail comparison.

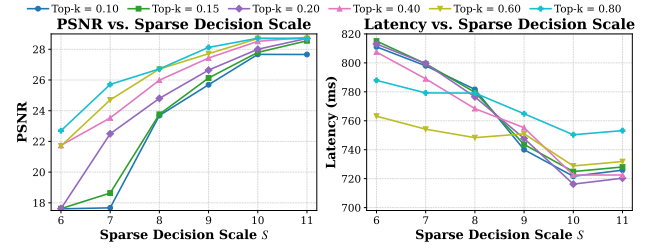


Figure 6. Ablation on sparse decision scale and Top-K.

Table 4. Ablation study on the attention sink and window size configuration in CSLA.

Attn. Sink	Wind. Size	Latency (ms) \downarrow	Sparsity \uparrow	GenEval \uparrow	PSNR \uparrow
Scale ≤ 5	[1, 3, 5]	0.4657	86.24%	0.727	25.657
	[3, 3, 3]	0.4480	86.77%	0.732	25.837
	[3, 5, 7]	0.5346	83.46%	0.730	25.897
	[5, 5, 5]	0.5266	83.77%	0.727	25.895
	[5, 7, 9]	0.6186	80.72%	0.729	25.944
	[7, 7, 7]	0.5960	81.18%	0.729	26.028
Scale ≤ 6	[7, 9, 11]	0.6889	78.06%	0.730	26.068
	[3, 5, 7]	0.6114	81.03%	0.733	26.928
Scale ≤ 7	[3, 5, 7]	0.6717	78.60%	0.730	27.174
Scale ≤ 8	[3, 5, 7]	0.7729	75.21%	0.733	27.582
w/o Sink	[3, 5, 7]	0.5008	84.68%	0.714	23.541

one dense sink block, shown as Fig. 4(d)) as attention sinks is crucial. Without early-scale global structural information, sparse attention at later large scales will lead to a degradation in generation quality, causing PSNR to drop from **25.897** to **23.541**. For the last 3 scales, smaller windows yield the highest sparsity and fastest kernel latency, but at the cost of degraded detail generation. Balancing quality and efficiency, we adopt scale ≤ 5 as the default sinks and use window size [3, 5, 7] for the last 3 scales.

Effectiveness of Cross-Scale Sparse Attention. We conduct an ablation study to analyze the effectiveness of our proposed plug-and-play efficient sparse attention modules. The results are shown in Tab. 5, introducing CS^4A alone yields a substantial acceleration of $1.46\times$ while maintaining competitive generation quality. A key design choice in CS^4A is whether to incorporate the cached output $\mathbf{O}_{\text{cache}}^{(k)}$ from the sparse decision scale to supplement subsequent large-scale sparse outputs. Experimental results demon-

Table 5. Ablation study on our **plug-and-play efficient sparse attention modules**.

Methods	Speedup \uparrow	Latency \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Infinity-8B [12]	1.00 \times	1.65s	-	-	-
+ CS^4A (w/o $\mathbf{O}_{\text{cache}}^{(k)}$)	1.46 \times	1.13s	25.665	0.837	0.131
+ CS^4A (w/ $\mathbf{O}_{\text{cache}}^{(k)}$)	1.43 \times	1.15s	26.359	0.860	0.114
+ $CLSA$ (block-wise)	1.57\times	1.05s	29.481	0.920	0.073

strate that adding $\mathbf{O}_{\text{cache}}^{(k)}$ leads to higher visual fidelity—improving PSNR by approximately 0.7—while incurring negligible latency overhead. Finally, integrating the CLSA further enhances spatial coherence, achieving the best overall performance with a **1.57 \times** speedup and superior generation quality. These results highlight the importance of intra- and inter-scale spatial locality in VAR attention, suggesting that leveraging such local priors provides a powerful basis for training-free acceleration.

5. Conclusion

We introduce SparVAR, a training-free acceleration framework for Visual AutoRegressive (VAR) models that exploits the intrinsic sparsity of cross-scale attention. By systematically analyzing pretrained VAR, we identify three key properties—*attention sinks*, *cross-scale activation similarity*, and *spatial locality*—and use them to design two plug-and-play sparse attention modules: Cross-Scale Self-Similar Sparse Attention (CS^4A), which predicts high-resolution scales sparse patterns from a sparse decision scale, and Cross-Scale Local Sparse Attention (CLSA), which enforces locality via a block-wise sparse kernel whose forward is over $5\times$ faster than FlashAttention. On 1024×1024 text-to-image generation with Infinity-8B, SparVAR achieves a **1.57 \times** speedup without skipping any scales while preserving almost all high-frequency details and human preference scores comparable to or better than the baseline. When combined with existing scale-skipping strategies, SparVAR further attains up to $2.28\times$ acceleration with competitive GenEval score and low-level metrics. These results demonstrate that leveraging cross-scale attention sparsity is a principled and effective direction to scale VAR inference efficiently.

References

- [1] Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv e-prints*, pages arXiv–2506, 2025. 3, 2
- [2] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion

- transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 2
- [3] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In *European Conference on Computer Vision*, pages 74–91. Springer, 2024. 2, 3
- [4] Pengtao Chen, Mingzhu Shen, Peng Ye, Jianjian Cao, Chongjun Tu, Christos-Savvas Bouganis, Yiren Zhao, and Tao Chen. δ -dit: A training-free acceleration method tailored for diffusion transformers. *arXiv preprint arXiv:2406.01125*, 2024. 4
- [5] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023. 3, 5
- [6] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022. 3
- [7] Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex attention: A programming model for generating optimized attention kernels. *arXiv preprint arXiv:2412.05496*, 2024. 5, 6
- [8] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. 2, 3
- [9] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 3
- [10] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023. 2, 7
- [11] Hang Guo, Yawei Li, Taolin Zhang, Jiangshan Wang, Tao Dai, Shu-Tao Xia, and Luca Benini. Fastvar: Linear visual autoregressive modeling via cached token pruning. *arXiv preprint arXiv:2503.23367*, 2025. 2, 4, 7, 8, 5, 6
- [12] 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 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15733–15744, 2025. 1, 2, 3, 4, 6, 7, 8, 9, 5
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 2
- [14] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020. 3

- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2
- [16] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 7, 4
- [17] Ting Jiang, Yixiao Wang, Hancheng Ye, Zishan Shao, Jingwei Sun, Jingyang Zhang, Zekai Chen, Jianyi Zhang, Yiran Chen, and Hai Li. Sada: Stability-guided adaptive diffusion acceleration. *arXiv preprint arXiv:2507.17135*, 2025. 4
- [18] Kumara Kahatapitiya, Haozhe Liu, Sen He, Ding Liu, Menglin Jia, Chenyang Zhang, Michael S Ryoo, and Tian Xie. Adaptive caching for faster video generation with diffusion transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15240–15252, 2025. 4
- [19] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 3
- [20] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 2
- [21] Jiajun Li, Yue Ma, Xinyu Zhang, Qingyan Wei, Songhua Liu, and Linfeng Zhang. Skipvar: Accelerating visual autoregressive modeling via adaptive frequency-aware skipping. *arXiv preprint arXiv:2506.08908*, 2025. 2, 4, 7, 8, 5
- [22] Kunjun Li, Zigeng Chen, Cheng-Yen Yang, and Jenq-Neng Hwang. Memory-efficient visual autoregressive modeling with scale-aware kv cache compression. *arXiv preprint arXiv:2505.19602*, 2025. 2, 4, 7, 8, 5
- [23] Lijiang Li, Huixia Li, Xiawu Zheng, Jie Wu, Xuefeng Xiao, Rui Wang, Min Zheng, Xin Pan, Fei Chao, and Rongrong Ji. Autodiffusion: Training-free optimization of time steps and architectures for automated diffusion model acceleration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7105–7114, 2023. 4
- [24] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:56424–56445, 2024. 2
- [25] Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snap-fusion: Text-to-image diffusion model on mobile devices within two seconds. *Advances in Neural Information Processing Systems*, 36:20662–20678, 2023. 4
- [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, 3
- [27] Feng Liu, Shiwei Zhang, Xiaofeng Wang, Yujie Wei, Haonan Qiu, Yuzhong Zhao, Yingya Zhang, Qixiang Ye, and Fang Wan. Timestep embedding tells: It’s time to cache for video diffusion model. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7353–7363, 2025. 4
- [28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, pages 1–22, 2025. 4
- [29] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for free. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15762–15772, 2024. 4
- [30] Ziqi Pang, Tianyuan Zhang, Fujun Luan, Yunze Man, Hao Tan, Kai Zhang, William T Freeman, and Yu-Xiong Wang. Randar: Decoder-only autoregressive visual generation in random orders. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 45–55, 2025. 2
- [31] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2
- [32] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 3
- [33] Ziran Qin, Youru Lv, Mingbao Lin, Zeren Zhang, Danping Zou, and Weiyao Lin. Head-aware kv cache compression for efficient visual autoregressive modeling. *arXiv preprint arXiv:2504.09261*, 2025. 4
- [34] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 3
- [35] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in Neural Information Processing Systems*, 32, 2019. 3
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 2, 3
- [37] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017. 3
- [38] Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. *Advances in Neural Information Processing Systems*, 37: 68658–68685, 2024. 5
- [39] Austin Silveria, Soham V Govande, and Daniel Y Fu. Chipmunk: Training-free acceleration of diffusion transformers with dynamic column-sparse deltas. In *ES-FoMo III: 3rd Workshop on Efficient Systems for Foundation Models*, 2025. 4, 5, 1
- [40] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 2, 3

- [41] Haotian Tang, Yecheng Wu, Shang Yang, Enze Xie, Junsong Chen, Junyu Chen, Zhuoyang Zhang, Han Cai, Yao Lu, and Song Han. HART: Efficient visual generation with hybrid autoregressive transformer. In *The Thirteenth International Conference on Learning Representations*, 2025. 2, 3, 6, 7
- [42] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems*, 37:84839–84865, 2024. 2, 3, 4
- [43] Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19, 2019. 5
- [44] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016. 3
- [45] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1747–1756. PMLR, 2016. 3
- [46] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in Neural Information Processing Systems*, 30, 2017. 3
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 3
- [48] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 2
- [49] Lei Wang, Yu Cheng, Yining Shi, Zhengju Tang, Zhiwen Mo, Wenhao Xie, Lingxiao Ma, Yuqing Xia, Jilong Xue, Fan Yang, et al. Tilelang: A composable tiled programming model for ai systems. *arXiv preprint arXiv:2504.17577*, 2025. 5
- [50] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 2
- [51] Yuqing Wang, Shuhuai Ren, Zhijie Lin, Yujin Han, Haoyuan Guo, Zhenheng Yang, Difan Zou, Jiashi Feng, and Xihui Liu. Parallelized autoregressive visual generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12955–12965, 2025. 2
- [52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 2
- [53] Felix Wimbauer, Bichen Wu, Edgar Schoenfeld, Xiaoliang Dai, Ji Hou, Zijian He, Artsiom Sanakoyeu, Peizhao Zhang, Sam Tsai, Jonas Kohler, et al. Cache me if you can: Accelerating diffusion models through block caching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6211–6220, 2024. 4
- [54] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 7, 8, 4, 6
- [55] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. In *International Conference on Representation Learning*, pages 21875–21895, 2024. 3
- [56] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 2
- [57] Jing Xiong, Gongye Liu, Lun Huang, Chengyue Wu, Taiqiang Wu, Yao Mu, Yuan Yao, Hui Shen, Zhongwei Wan, Jinfa Huang, et al. Autoregressive models in vision: A survey. *arXiv preprint arXiv:2411.05902*, 2024. 2
- [58] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023. 7, 8
- [59] Evelyn Zhang, Bang Xiao, Jiayi Tang, Qianli Ma, Chang Zou, Xuefei Ning, Xuming Hu, and Linfeng Zhang. Token pruning for caching better: 9 times acceleration on stable diffusion for free. *arXiv preprint arXiv:2501.00375*, 2024. 4
- [60] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 2
- [61] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. *arXiv preprint arXiv:2412.20404*, 2024. 2
- [62] Chang Zou, Evelyn Zhang, Runlin Guo, Haohang Xu, Conghui He, Xuming Hu, and Linfeng Zhang. Accelerating diffusion transformers with dual feature caching. *arXiv preprint arXiv:2412.18911*, 2024. 4

SparVAR: Exploring Sparsity in Visual AutoRegressive Modeling for Training-Free Acceleration

Supplementary Material

A. Algorithmic Details of CS^4A

A.1. Hardware-Efficient Sparse Computation

As illustrated in Fig. 7, the pipeline of CS^4A comprises two distinct phases: sparse pattern identification and hardware-accelerated sparse computation. At the sparse decision scale S , we perform full dense attention to capture the global dependency structure. To avoid the computational overhead associated with sorting attention scores for individual queries, we adopt a block-wise Column-Sum strategy [39]. We partition the attention map $\mathbf{P}^{(S)}$ into contiguous query blocks of size C , as highlighted by the red dashed boxes in Fig. 7. Within each query block, we aggregate attention scores along the query dimension to derive a cumulative importance score for each key column. Subsequently, the sparse indices $\text{inds}^{(S)}$ are determined by performing a Top- K selection on these column sums. These indices are then projected to the target scale K via our sparse index mapping strategy $\mathcal{M}_{S \rightarrow K}(\cdot)$, enabling efficient sparse computation at high-resolution scales.

A primary challenge in sparse attention lies in the inefficiency of unstructured sparsity, where random sparse matrix multiplication fails to effectively leverage the Tensor Cores on modern GPUs. Furthermore, runtime overheads such as dynamic mask computation and cache maintenance can introduce significant latency. To address this, we adopt the Tile Packing technique [39], which maps sparse computation into compacted dense computation. We extend their efficient custom kernels to support our causal cross-scale attention mechanism. Specifically, guided by the mapped sparse indices $\text{inds}^{(K)}$, the kernel **gathers** discrete, non-contiguous key and value vectors from the GPU global memory (HBM). These vectors are then **packed** into contiguous dense tiles within the GPU shared memory (SRAM). Once resident in SRAM, these packed tiles constitute a logically dense matrix (as shown in the bottom-left of Fig. 7). This transformation enables the GPU to execute standard dense General Matrix Multiplications (GEMMs) using Tensor Cores, thereby achieving peak computational throughput.

A.2. Cross-Scale Sparse Index Mapping

To efficiently propagate the sparse activation patterns from the decision scale S to a higher-resolution target scale K , we propose a two-stage mapping mechanism, denoted as $\mathcal{M}_{S \rightarrow K}(\cdot)$. This mapping addresses two fundamental challenges: the quadratic expansion of the query grid and the

scale drift of historical key-value pairs. Specifically, we decompose \mathcal{M} into two coupled transformations: *Query Block Homography* and *Relative Scale Alignment*.

Query Block Homography. Following the hardware-efficient Column-Sum design [39], query tokens are partitioned into contiguous blocks of size C (e.g., $C = 192$). As the spatial resolution expands from $h_S \times w_S$ to $h_K \times w_K$, the number of query blocks increases from $G_S = \lceil N_S/C \rceil$ to $G_K = \lceil N_K/C \rceil$. Let $\mathcal{G}_S = \{0, \dots, G_S - 1\}$ and $\mathcal{G}_K = \{0, \dots, G_K - 1\}$ denote the sets of query block indices at scales S and K , respectively. We model the mapping from a target block $g_K \in \mathcal{G}_K$ to a source block $g_S \in \mathcal{G}_S$ as a nearest-neighbor interpolation within the normalized coordinate space $[0, 1]$. Formally, the mapping is defined as:

$$\phi(g_K) = \left\lfloor \frac{g_K + 0.5}{G_K} \cdot G_S - 0.5 \right\rfloor_{\mathcal{G}_S},$$

where $\lfloor \cdot \rfloor_{\mathcal{G}_S}$ denotes rounding to the nearest integer clipped within the bounds of \mathcal{G}_S . This transformation ensures that query tokens at the finer scale k inherit the sparsity pattern from their spatially corresponding region at the coarser scale S , thereby preserving the global attention structure.

Relative Scale Alignment for KVs. The KV cache in VAR is a concatenation of tokens from all historical scales, i.e., $\mathbf{K}_{\leq S} = \text{Concat}(\mathbf{K}_1, \dots, \mathbf{K}_S)$. A flattened global KV index j is structurally ambiguous without context. We propose a *Decompose-Align-Project* mechanism to map these indices:

- **Decomposition.** We first decompose a global index j into its constituent scale l and local spatial offset δ . Let $\{C_i\}_{i=1}^S$ be the cumulative token counts up to each scale (e.g., $C_3 = N_1 + N_2 + N_3 = 1 \times 1 + 2 \times 2 + 4 \times 4 = 21$). The source scale l and offset δ_S are retrieved via:

$$l = \max\{i \mid C_i \leq j\}, \quad \delta_S = j - C_l. \quad (4)$$

- **Relative Alignment.** Based on the cross-scale self-similarity observation, the sparse attention pattern depends on the *relative* depth between the query and key scales. For a query at scale K , attending to the m -th previous scale should mirror the behavior of scale S attending to its m -th previous scale. Thus, the target scale level l' at step K is determined by preserving the relative scale offset:

$$l' = K - (S - l).$$

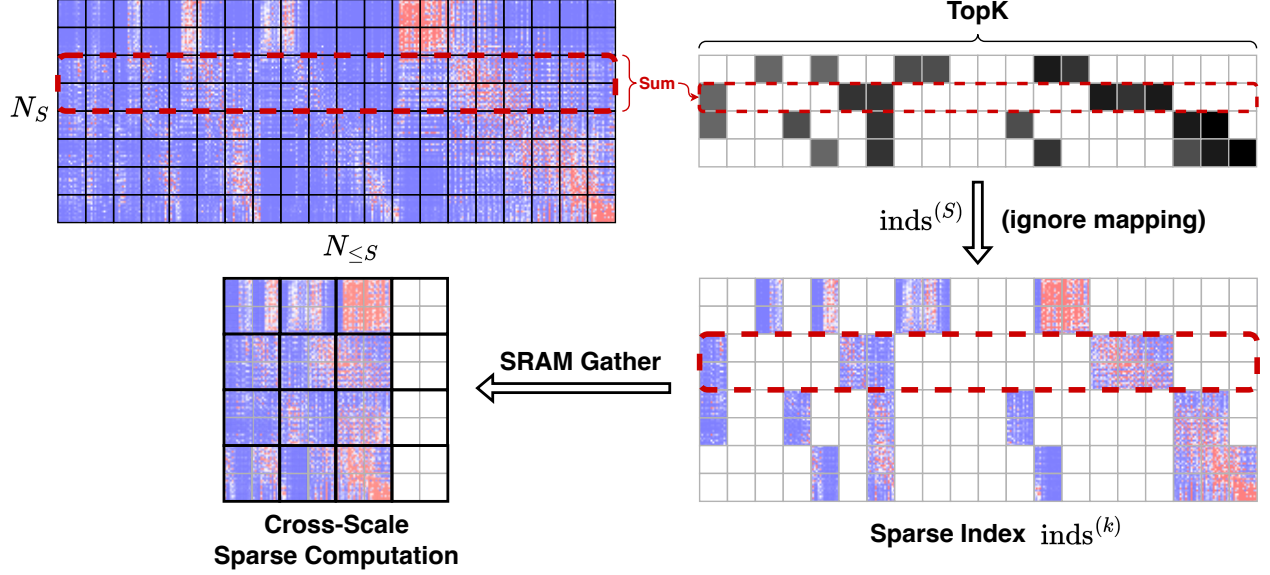


Figure 7. The illustration of CS^4A . The sparse index mapping process is omitted for clarity.

- **Spatial Projection.** Finally, we project the local spatial offset δ_S from the resolution of scale l ($h_l \times w_l$) to the resolution of the target scale l' ($h_{l'} \times w_{l'}$). To preserve 2D spatial locality, we de-linearize δ_S into 2D coordinates (u, v) , perform bilinear scaling, and re-linearize:

$$u' = \lfloor u \cdot \frac{h_{l'}}{h_l} \rfloor, \quad v' = \lfloor v \cdot \frac{w_{l'}}{w_l} \rfloor, \quad \delta_K = u' \cdot w_{l'} + v'. \quad (5)$$

The final mapped global index is reconstructed as $j' = C_{l'} + \delta_K$.

Attention Sink. To mitigate the loss of critical global context during sparse selection, we explicitly enforce a **sink token preservation** strategy. Let $\mathcal{A}_{\text{sink}} = \{0, \dots, N_{\text{sink}} - 1\}$ denote the set of indices corresponding to tokens from the initial scales. The final sparse index set for the target scale K is constructed as the union of the mapped dynamic indices and these static sink indices:

$$\text{inds}^{(K)} = \mathcal{A}_{\text{sink}} \cup \mathcal{M}_{S \rightarrow K}(\text{inds}^{(S)}).$$

This ensures that the accelerated sparse attention mechanism retains a stable attention sink, thereby preventing the collapse of semantic coherence during the generation process.

Effectiveness Analysis. To validate the effectiveness of the proposed sparse index mapping strategy, we visualize the alignment between the predicted sparse patterns and the full attention activations of the baseline model in Fig. 8. Specifically, we visualize the sparse indices derived from the sparse decision scale S (highlighted by red dashed

boxes) alongside the mapped indices at subsequent high-resolution scales, allowing for a direct examination of the sparse patterns predicted by CS^4A . As illustrated across various layers and heads, the mapped sparse indices (depicted in gold) exhibit a remarkable spatial correspondence with the ground-truth activation hotspots (red regions) of the baseline model. Notably, this alignment persists even in the final scales where the resolution expands significantly, confirming that the structural self-similarity of cross-scale attention is robustly preserved by our mapping strategy. This coverage of salient attention activation regions ensures that the sparse computation approximates the full attention mechanism with minimal error. It fundamentally explains why CS^4A achieves substantial acceleration while maintaining the generation of high-frequency details, as the model continues to attend to the most critical contextual information during the fine-grained refinement stages.

A.3. Difference from Chipmunk.

Chipmunk achieves training-free acceleration by exploiting the temporal redundancy and the similarity of attention activations across denoising steps inherent in Diffusion Transformers (DiTs) [1, 20, 48]. However, the premise that latent variables maintain a constant spatial resolution during denoising does not hold for the unique next-scale prediction paradigm of VAR. Fundamentally, VAR is characterized by a progressive increase in token count as spatial resolution grows during inference, standing in stark contrast to DiT, where the sequence length remains invariant. The motivation behind our proposed CS^4A lies in the cross-scale self-similarity intrinsic to VAR models, namely that attention activation patterns at smaller scales exhibit similarities

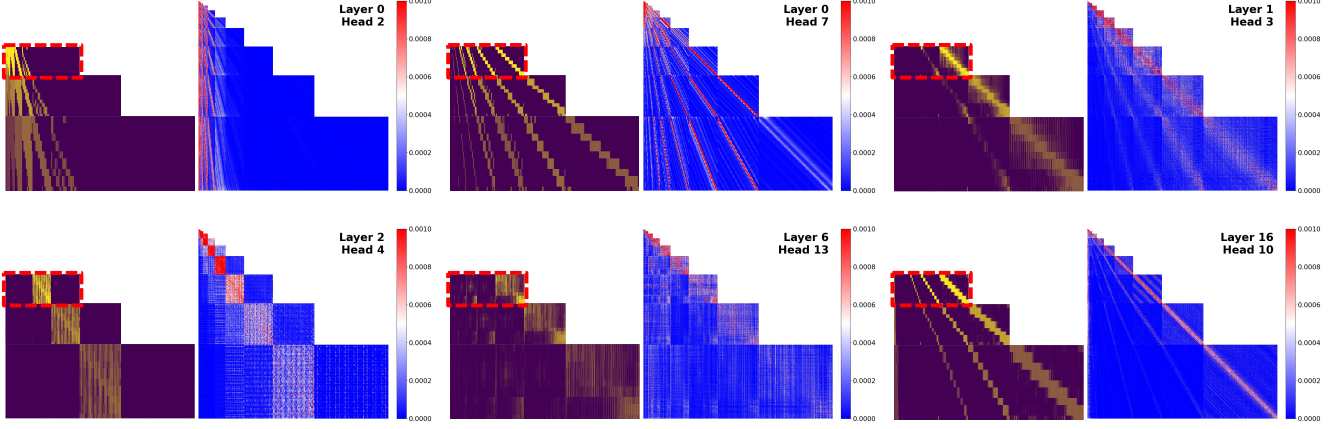


Figure 8. Visualization of cross-scale sparse index mapping. The red dashed box highlights the sparse indices identified at the sparse decision scale S , while the blocks below display the corresponding sparse indices projected onto subsequent high-resolution scales. Evidently, the mapped sparse indices precisely cover the most salient attention activation regions of the baseline model.

to those at larger scales in corresponding regions. To the best of our knowledge, SparVAR is the first work to identify and exploit this attention activation pattern redundancy for accelerating visual autoregressive generation, exploring a dimension of sparsity orthogonal to the temporal sparsity found in DiTs. Furthermore, to effectively propagate this sparsity across scales, we propose an efficient sparse index mapping strategy that projects sparse patterns from low-resolution to high-resolution scales via rigorous geometric transformations. From an engineering perspective, we extend Chipmunk’s highly optimized sparse kernels—originally designed solely for standard self-attention—to support the causal cross-scale attention required by VAR. This generalization unlocks the potential of sparse computing for a broader class of autoregressive models beyond the fixed-length setting of DiTs.

B. Complexity of Cross-Scale Sparse Attention

We analyze the time complexity of the proposed CS^4A and $CSLA$, benchmarking them against standard cross-scale full attention. Let N_k denote the number of tokens at the current scale k , and $N_{\leq k}$ represent the cumulative number of tokens across all preceding scales. The feature dimension is denoted by d . For a given scale k , following Eq. 1, the computational complexity of full dense attention is given by:

$$\mathcal{O}(dN_kN_{\leq k}).$$

CS^4A . CS^4A initially performs full dense attention at the sparse decision scale S to identify salient sparse patterns. The complexity of full attention at scale S is $\mathcal{O}(dN_SN_{\leq S})$, consistent with the baseline model. Since S typically corresponds to a mid-resolution scale, this computational overhead is relatively marginal. Additionally, computing the column-sum tensor $D^{(S)}$ per query block

introduces a complexity of $\mathcal{O}(dG_SN_{\leq S})$, where $G_S = \lceil N_S/C \rceil$. Consequently, the total complexity at the sparse decision scale S is given by:

$$\mathcal{O}^{(S)} = \mathcal{O}(dN_SN_{\leq S}) + \mathcal{O}(dG_SN_{\leq S}).$$

In practice, hardware-friendly block size C (e.g., 128 or 192) makes the second term negligible compared to the first.

For subsequent high-resolution scales $k > S$, we employ cross-scale sparse computation, which restricts attention calculation to the Top-K selected KV pairs. The process begins with sparse index mapping. Mapping the query blocks and KV indices incurs a time complexity of $\mathcal{O}(N_k)$, as it involves only lightweight coordinate transformations and gather operations. Let α denote the sparsity ratio (e.g., retaining the top 20% of KVs). For any given query, the number of active KVs is $\alpha \cdot N_{\leq k} \ll N_{\leq k}$. Therefore, the complexity of the sparse computation is reduced to

$$\mathcal{O}(d\alpha N_kN_{\leq k}).$$

This implies a linear relationship between total complexity and the sparsity ratio. For small α , CS^4A achieves significant acceleration. Theoretically, as N_k increases, the complexity reduction factor asymptotically approaches $1/\alpha$.

$CSLA$. $CSLA$ decomposes all historical KV caches into a dense attention sink and a block-wise local sparse region. The first N_{sink} KV tokens from early scales are always attended with full dense attention, preserving their role as global structural anchors [55]. The remaining $N_{\leq k} - N_{\text{sink}}$ KV tokens are arranged as 2D grids per historical scale, where each scale h is equipped with a scale-dependent local window size w_h (with radius $r_h = \lfloor w_h/2 \rfloor$). Consequently, any query at scale k interacts with at most

$$N_{\text{local}} \leq \sum_{h \leq k} w_h^2$$

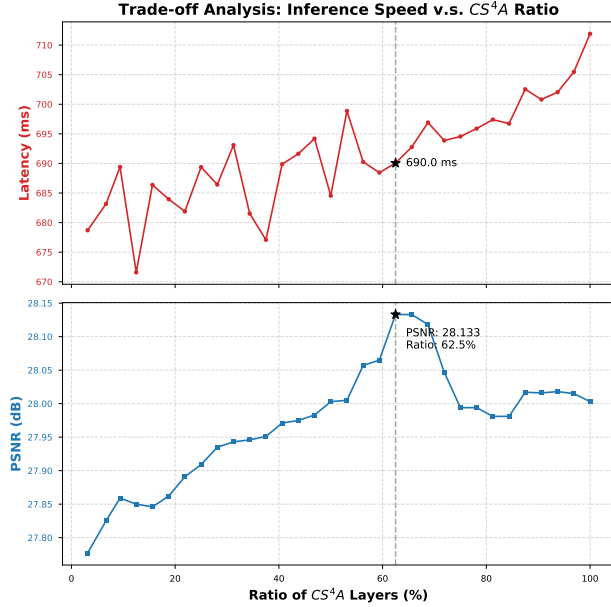


Figure 9. Configuration trade-off between CS^4A and $CSLA$.

KV tokens in the sparse region. In our design, the window sizes, the window sizes w_h are kept small—for example, the last three scales use $\{3, 5, 7\}$ —so N_{local} is effectively a constant upper bound. The computational complexity of $CSLA$ at scale k therefore becomes

$$\mathcal{O}(dN_k(N_{\text{sink}} + N_{\text{local}})).$$

Moreover, since N_{sink} covers only a few early scales (e.g. the first 5 scales contain just 121 KV tokens), we have $N_{\text{sink}} + N_{\text{local}} \ll N_{\leq k}$ at late high-resolution scales.

C. Implementation Details.

The baseline VAR model operates across a total of 13 resolution scales, corresponding to 13 iterations of next-scale autoregressive prediction. We configure the sparsity settings as follows: (i) In the setting **without scale skipping**, we designate Scale 11 as the sparse decision scale and apply sparsification to the last 2 scales to accelerate inference. (ii) When **skipping the last two scales**, Scale 10 serves as the sparse decision scale, with sparsification applied exclusively to Scale 11. (iii) For the **sample-adaptive skipping strategy**, we similarly adopt Scale 10 as the sparse decision scale, sparsifying any subsequent scales that are not skipped. Regarding the configuration of sparse attention modules, the application ratio of CS^4A to $CSLA$ across all attention layers is set to 6:4. Finally, for the acceleration performance analysis, to ensure a fair comparison with the baseline model, we measure throughput using a batch size of 1. To facilitate reproducibility, we will release our

efficient kernels and full implementation code upon acceptance.

D. Additional Experiments

Comparison on DPG-Bench. To further validate the effectiveness of SparVAR, we provide additional quantitative results on the DPG-Bench [16], which evaluates prompt-following capabilities across multiple dimensions. The results are summarized in Tab. 6. For the 2B model, SparVAR achieves a competitive overall score of **82.858** without skipping scales, surpassing the baseline and FastVAR while delivering a superior **1.38×** speedup. When combined with scale-skipping strategies, our method maintains robust performance, achieving the highest score of **82.891** under the dynamic skipping setting with a **1.58×** speedup, demonstrating an effective balance between acceleration and semantic fidelity. For the larger 8B model, SparVAR achieves an overall score of **86.411** without scale skipping, comparable to the baseline (86.440) and outperforming FastVAR and ScaleKV, all while providing a substantial **1.57×** acceleration. Notably, when integrated with scale-skipping strategies, SparVAR consistently outperforms other baselines. In the “skip last 2 scales” setting, it attains the highest score of **86.468**—exceeding even the baseline—alongside a remarkable speedup. The above results demonstrate the generality and effectiveness of our method.

Qualitative Comparison. We compare SparVAR against the baseline Infinity-2B and other training-free acceleration methods across complex generation scenarios from the HPSv2.1 [54] benchmark. As shown in Fig. 10, in the setting without scale skipping, SparVAR achieves a **1.38×** speedup while generating images with visual fidelity nearly identical to the baseline model. This is particularly evident in the second row, where our method accurately reconstructs the fine-grained square window panes of the white building—a structural detail that both FastVAR and ScaleKV fail to preserve, resulting in blurred or distorted features. Crucially, the robustness of SparVAR extends to the aggressive skip-scale setting. When skipping the last two high-resolution scales to achieve a **1.70×** speedup, SparVAR still successfully generates coherent structural details (e.g., the square windows), benefiting from the reinforced spatial locality in our sparse attention mechanism. In contrast, FastVAR exhibits degradation, introducing severe artifacts and failing to maintain the geometric integrity of the scene. These results qualitatively confirm that SparVAR’s cross-scale sparse attention is highly effective at compressing redundant computation while selectively preserving the critical high-frequency information necessary for photorealistic generation.

Table 6. Quantitative comparison on the DPG-Bench dataset. The best results in each setting are highlighted in **bold**.

Methods	Acceleration		DPG-Bench					
	#Scales ↓	Speedup ↑	Global ↑	Entity ↑	Attribute ↑	Relation ↑	Other ↑	Overall ↑
<i>w/o skip scales</i>								
Infinity-2B [12]	13	1.00×	88.499	88.120	88.957	88.767	84.521	82.839
FastVAR [11]	13	1.01×	77.291	89.168	88.537	91.932	84.637	82.723
ScaleKV [22]	13	0.70×	78.088	88.550	87.495	90.229	89.468	82.885
SparVAR (Ours)	13	1.38×	89.021	89.174	87.274	91.015	84.405	82.858
<i>w/ skip last 2 scales</i>								
FastVAR [11]	11	1.33×	81.423	88.267	88.184	90.127	89.351	82.651
SparVAR (Ours)	11	1.70×	83.864	88.945	88.274	90.883	81.888	82.663
<i>w/ dynamic skip last scales</i>								
SkipVAR-2B [21]	10~13	1.33×	85.273	89.115	89.289	89.519	84.836	82.877
SparVAR (Ours)	10~13	1.58×	89.216	88.896	88.004	89.279	89.366	82.891
<i>w/o skip scales</i>								
Infinity-8B [12]	13	1.00×	92.144	89.741	90.376	92.084	91.862	86.440
FastVAR [11]	13	1.14×	81.784	90.195	89.311	94.084	84.496	86.343
ScaleKV [22]	13	0.67×	81.876	90.041	89.666	94.464	92.675	86.333
SparVAR (Ours)	13	1.57×	87.306	91.510	90.862	91.565	90.193	86.411
<i>w/ skip last 2 scales</i>								
FastVAR [11]	11	1.79×	87.295	91.466	90.048	91.113	90.308	86.332
SparVAR (Ours)	11	2.28×	91.729	91.068	91.031	90.354	89.072	86.468
<i>w/ dynamic skip last scales</i>								
SkipVAR-8B [21]	10~13	1.71×	92.648	90.591	91.032	90.088	86.008	86.293
SparVAR (Ours)	10~13	2.05×	88.227	91.863	90.703	92.239	86.758	86.409

Table 7. Quantitative comparison of memory consumption. ‘‘Torch Mem.’’ refers to peak PyTorch allocated memory, and ‘‘NV Mem.’’ refers to NVIDIA physical memory. The best results in each setting are highlighted in **bold**.

Methods	2B model				8B model			
	GenEval	Speedup ↑	Torch Mem. ↓	NV Mem. ↓	GenEval ↑	Speedup ↑	Torch Mem. ↓	NV Mem. ↓
<i>w/o skip scales</i>								
Infinity [12]	0.736	1.00×	15.78GB	28.63GB	0.798	1.00×	37.11GB	55.75GB
FastVAR [11]	0.712	1.01×	15.91GB	28.45GB	0.792	1.14×	37.41GB	54.29GB
ScaleKV [22]	0.732	0.70×	10.98GB	17.01GB	0.793	0.67×	22.71GB	29.64GB
SparVAR (Ours)	0.738	1.38×	16.02GB	23.21GB	0.796	1.57×	37.48GB	45.45GB
+ KV comp.	0.730	1.29×	12.74GB	18.51GB	0.793	1.31×	30.42GB	35.18GB
<i>w/ skip last 2 scales</i>								
FastVAR [11]	0.716	1.33×	11.65GB	19.73GB	0.79	1.79×	26.19GB	37.31GB
SparVAR (Ours)	0.725	1.70×	11.12GB	17.04GB	0.800	2.28×	26.42GB	29.73GB
<i>w/ dynamic skip last scales</i>								
SkipVAR [21]	0.730	1.33×	12.50GB	23.32GB	0.789	1.71×	30.05GB	44.00GB
SparVAR (Ours)	0.732	1.58×	12.50GB	16.98GB	0.796	2.05×	30.00GB	34.10GB

Memory Consumption. We also provide a comprehensive memory comparison for generating 1024×1024 images in the Tab. 7, including peak PyTorch active memory (Torch Mem.) and NVIDIA physical memory (NV Mem.). The memory reduction in FastVAR and SkipVAR rely heavily on scale-skipping strategies. Notably, FastVAR (w/o skip scales) exhibits memory usage nearly to the baseline, as the overhead of caching intermediate scale tokens for restoration offsets the benefits of token pruning. While ScaleKV achieves the lowest memory, its design cannot be accelerated by modern efficiency attention kernels (e.g., FlashAt-

tention), resulting in inference speeds significantly slower than the baseline. Our sparse attention is optimized primarily for **inference latency** rather than memory minimization. Some reduction in physical memory originates from the efficient memory access strategies of the kernels. As a ‘plug-and-play’ module, SparVAR is compatible with KV compression methods (e.g., retaining only sinks and local scales KV Cache, denoted as ‘+ KV Comp.’ in Tab. 7). We leave further dedicated memory optimizations to future work.



Figure 10. Qualitative comparison of complex scene generation on HPSv2.1 [54] benchmark. Zoom in for fine-detail visualization.

Table 8. Quantitative comparison on efficiency and quality on 1024×1024 image generation. Note that the efficiency of HART baseline is tested under FlashAttention.

Methods	Acceleration				GPU Memory		GenEval	DPG-Bench	HPSv2.1	ImageReward	Quality Metrics		
	#Scales ↓	Speedup ↑	Latency ↓	Throughput ↑	Torch Mem. ↓	NV Mem. ↓					PSNR ↑	SSIM ↑	LPIPS ↓
HART [41]	14	1.00×	446.30ms	2.24it/s	19.35	27.93	0.509	74.754	29.070	0.661	-	-	-
FastVAR [11]	14	1.11×	403.33ms	2.48it/s	19.33	25.72	0.504	74.762	27.680	0.602	20.850	0.704	0.316
SparVAR (Ours)	14	1.16×	383.11ms	2.61it/s	19.29	26.98	0.507	75.625	29.140	0.680	21.157	0.735	0.240

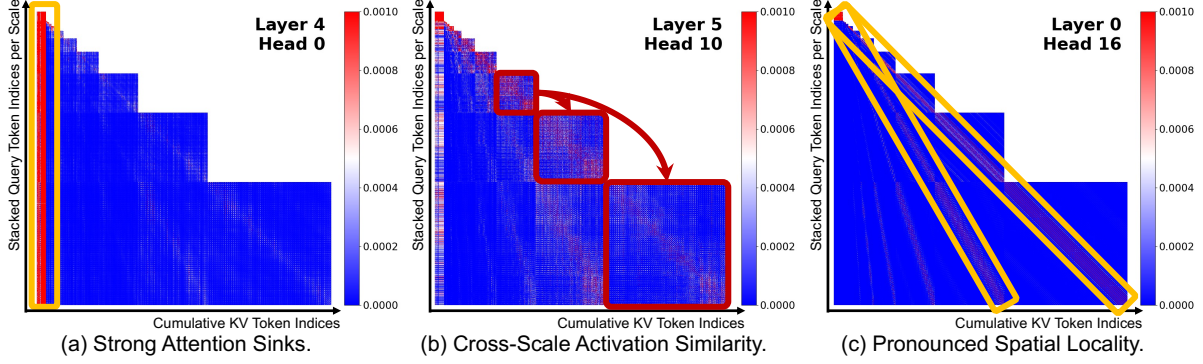


Figure 11. Visualization of attention activation patterns in the HART [41] across different layers and heads.

Generalizability of SparVAR. We further validated SparVAR on another representative VAR-based text-to-image model HART [41]. Our aim is to demonstrate that SparVAR is not custom-designed for a specific VAR architecture (e.g., Infinity [12]), but rather serves as a plug-and-play module for accelerating visual generation models based on the next-scale prediction paradigm. Fig. 11 illustrates the sparse activation patterns at each scale of the HART model, which exhibit three key properties similar to those observed in Infinity. It is worth noting that since the first scale in HART is not a 1×1 token map but comprises text tokens with a fixed sequence length of 300, its attention sink phenomenon is even more stronger compared to Infinity. Consequently, we adjust the configuration of the attention sink scales for our cross-scale sparse attention. Specifically, in addition to the first 5 scales, the sink region must consistently include the initial text tokens. Note that, following the setting of FastVAR, this model without scale skipping. Shown in Tab. 8, our SparVAR achieves speedup and consistently outperforms FastVAR across 4 high-level benchmarks, even exceeding the baseline. Superior low-level metrics further confirm robust detail preservation.

Optimal Configuration of Sparse Attention Modules.

As illustrated in Fig. 12, we further observe that in pre-trained VAR models, the attention activation pattern evolves with layer depth, shifting from a relatively diverse distribution to a pronounced local pattern concentrated along the cross-scale diagonal. Consequently, the dispersed activation patterns in shallower layers are better suited for CS^4A , which dynamically predicts sparse structures, whereas the fixed, strong locality in deeper layers is more effectively handled by the specialized CSLA. To determine the optimal configuration, we conduct a systematic ablation study on the layer-wise allocation of CS^4A and CSLA. Starting from a baseline where all attention layers utilize CSLA, we progressively substitute them with CS^4A from the shallowest to the deepest layers.

The quantitative results, illustrated in Fig. 9, reveal two key insights: **1) Synergistic effect on generation quality.** As the proportion of CS^4A gradually increases, we observe a corresponding steady improvement in PSNR. This trend aligns with our analysis, indicating that for the relatively dispersed sparse activation patterns in shallow layers, the cross-scale mapping of CS^4A achieves superior pattern prediction compared to the fixed local window sparse of CSLA. However, extending CS^4A to the deepest layers (beyond 70%) leads to a slight decline in PSNR. This suggests that deep layers are indeed specialized for high-frequency texture refinement, where the strong inductive bias of spatial locality enforced by CSLA is more effective than the mapped sparse patterns. **2) Latency-Quality trade-off.** In terms of efficiency, replacing CSLA with CS^4A introduces a marginal increase in latency, primarily due to the overhead of index mapping and memory gathering compared to the highly optimized block-wise kernel of CSLA. Nevertheless, the gain in visual fidelity (+0.35 dB PSNR) justifies this minimal cost. Consequently, we adopt a hybrid configuration with about 60% CS^4A and 40% CSLA as the default setting for 2B and 8B models, ensuring optimal performance.

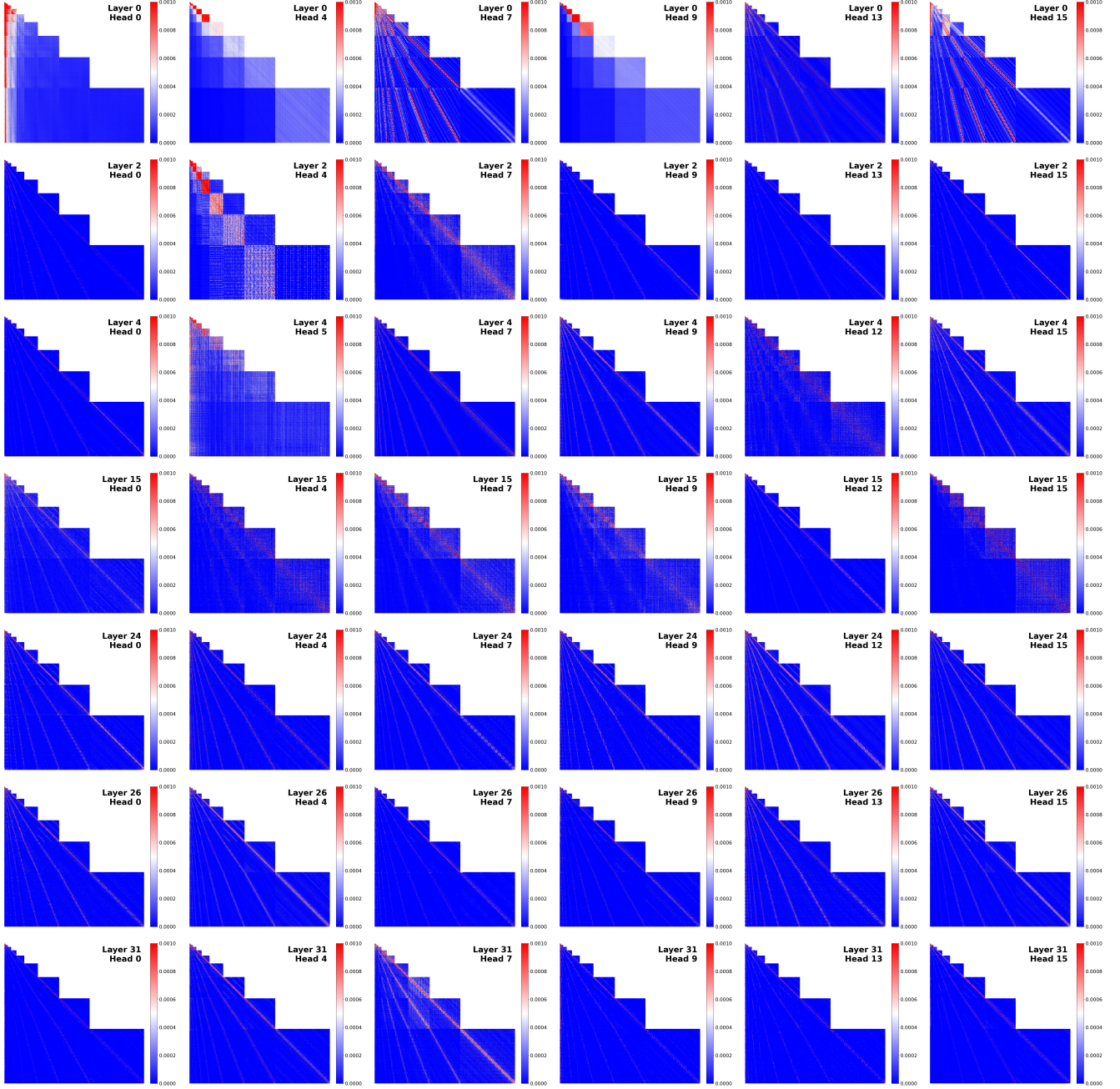


Figure 12. Visualization of the evolving attention activation patterns in VAR across layer depths.