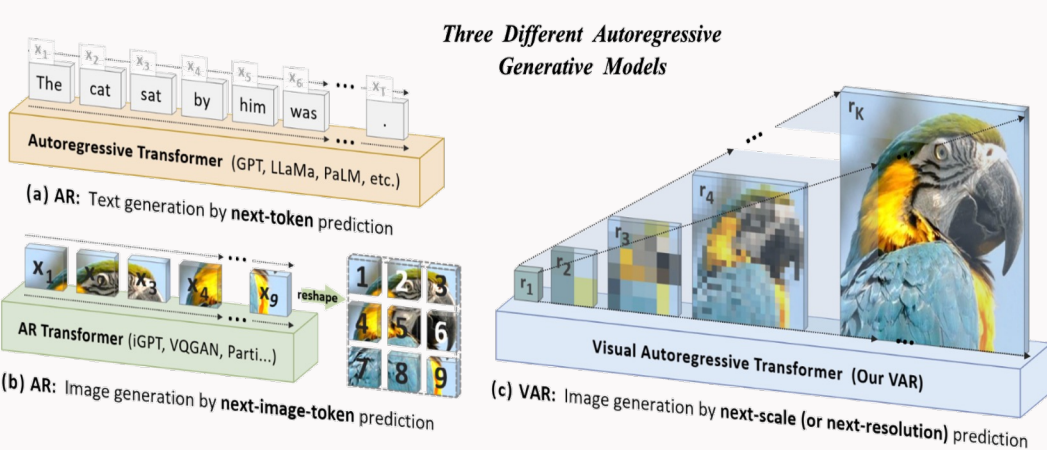


# ToProVAR: Efficient Visual Autoregressive Modeling via Tri-Dimensional Entropy-Aware Pruning

Jiayu Chen, Ruoyu Lin, Jingxin Li, Xiang Chen

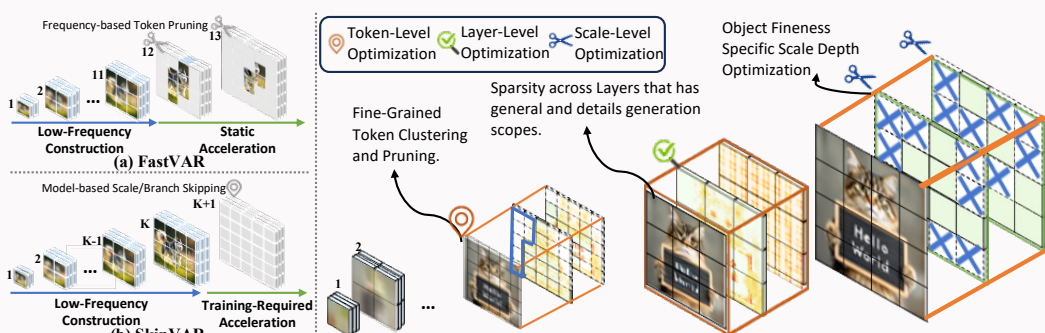
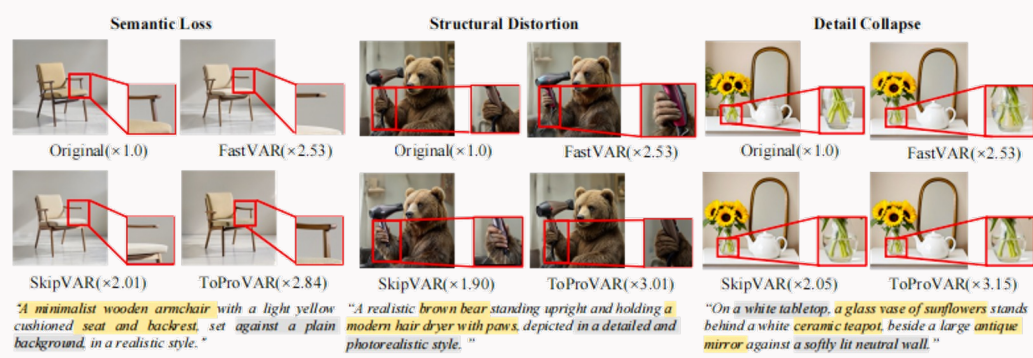
School of Computer Science, Peking University

## Background



- Traditional AR: **token-by-token**, strictly sequential → **slow & non-parallelizable**.
- VAR: generate **coarse** → **fine multi-scale tokens** → faster, but **cost grows exponentially** with scale.

## Problem



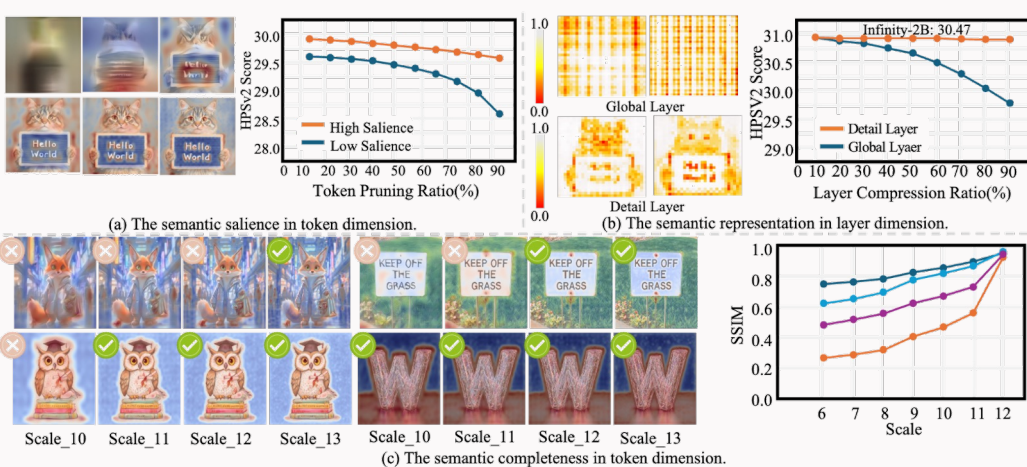
## Limitations of 1D Optimizations (FastVAR&SkipVAR)

- Only optimize one dimension (token / scale).
- Cannot capture **semantic relations** → structure distortion.
- Remove tokens that need **deeper scale** → detail collapse.

## 3D Sparsity Optimizations (ToProVAR)

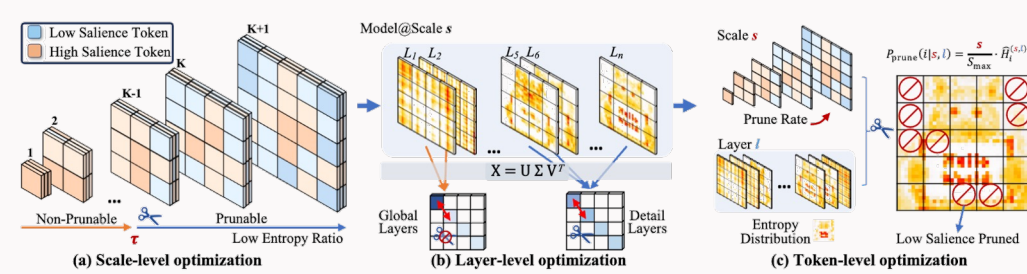
- We propose a **multi-dimensional optimization framework** that analyzes and prunes computation across **token, layer, and scale** dimensions.

## Analysis



- **Token-level:** Few tokens carry key semantics; pruning **90%** low-saliency tokens yields **<2%** quality loss.
- **Layer-level:** Global layers are pruning-sensitive, while detail layers can be pruned by **up to 90%**.
- **Scale-level:** Objects require different scale depths for correct semantic modeling.

## Method



## Scale-level optimization via Semantic Finess.

As scales grow, semantics stabilize.  
Measure stability using the **low-entropy ratio**:

$$\rho_s = \frac{|\{i | H_i^s < \bar{H}^s\}|}{N_s}$$

Pruning starts once

$$D = \min\{s | \rho_s \geq \tau\}$$

→ Only scales deeper than  $D$  are pruned.

## Layer-level optimization via Semantic Scope.

Use SVD to distinguish **Global Layers** vs. **Detail Layers**:

$$X = U\Sigma V^T, \quad \rho^{(l,s)} = \frac{\sigma_1}{\sigma_2}$$

Define representation score:

$$\mathcal{R}^{(l,s)} = \exp(-\beta(\rho^{(l,s)} - 1))$$

→ Prune only layers with **high  $\mathcal{R}^{(l,s)}$**  (Detail Layers).

## Token-level optimization via Fine-grained Semantic Saliency.

Normalize token entropy:

$$\hat{H}_i^{(l,s)} = \frac{H_i^{(l,s)}}{\sum_j H_j^{(l,s)}}$$

Compute unified pruning tendency:

$$q_i^{(l,s)} = \phi(s) \cdot \mathcal{R}^{(l,s)} \cdot \hat{H}_i^{(l,s)}, \quad \phi(s) = \frac{s}{S_{\max}}$$

Retention probability:

$$P_{\text{keep}}(i | s, l) = \begin{cases} 1, & s < D, \\ 1 - \text{clip}(\alpha_{\min} + (\alpha_{\max} - \alpha_{\min})q_i, 0, 1), & s \geq D. \end{cases}$$

→ High-entropy tokens in deep scales & detail layers are pruned first.

## Main Result

### Quantitative comparison on GenEval and DPG.

Methods	GenEval				DPG				Latency(s)↓	Speedup
	Two Obj.	Position	Color	Attri.	Overall↑	Entity	Relation	Attribute		
Infinity-2B	79.01	24.00	58.00	0.69	90.81	88.19	87.89	83.41	2.10	1.0 ×
+FastVAR	78.79	27.75	59.50	0.68	88.86	91.57	87.46	83.39	0.80	2.6 ×
+SkipVAR	76.77	26.50	57.50	0.67	89.30	87.07	87.01	82.94	1.10	2.0 ×
+ToProVAR	<b>78.80</b>	<b>29.50</b>	<b>62.00</b>	<b>0.69</b>	<b>87.39</b>	<b>88.92</b>	<b>90.01</b>	<b>83.07</b>	<b>0.61</b>	<b>3.4 ×</b>
Infinity-8B	96.97	61.00	75.00	0.83	90.92	93.57	88.83	86.68	4.86	1.0 ×
+FastVAR	94.19	57.00	75.25	0.81	90.80	92.30	90.40	86.50	2.01	2.4 ×
+SkipVAR	94.94	57.50	76.50	0.82	89.71	90.52	90.02	86.44	2.11	2.3 ×
+ToProVAR	<b>94.95</b>	<b>61.00</b>	<b>76.00</b>	<b>0.83</b>	<b>91.11</b>	<b>90.39</b>	<b>91.04</b>	<b>86.70</b>	<b>1.78</b>	<b>2.7 ×</b>

### Quantitative Comparison on HPSv2 and ImageReward.

Methods	HPSv2.1				ImageReward↑	Latency(s)↓	Speedup
	Photo	Concept-Art	Anime	Paintings			
Infinity-2B	29.40	30.38	31.72	30.39	30.47	0.94	1.57
+FastVAR	28.86	29.90	31.12	29.92	29.95	0.92	0.62
+SkipVAR	29.25	<b>30.25</b>	31.50	<b>30.45</b>	<b>30.39</b>	0.93	0.87
+ToProVAR	<b>29.26</b>	<b>30.15</b>	<b>31.44</b>	<b>30.23</b>	<b>30.27</b>	<b>0.93</b>	<b>0.58</b>
Infinity-8B	29.42	31.27	32.45	30.83	30.99	1.04	5.31
+FastVAR	<b>29.87</b>	30.42	31.80	29.89	30.24	1.02	1.97
+SkipVAR	29.09	30.86	32.04	<b>30.55</b>	<b>30.64</b>	1.03	2.65
+ToProVAR	29.19	<b>30.89</b>	<b>32.09</b>	30.24	30.58	<b>1.04</b>	<b>1.75</b>

### Quantitative comparisons of FID and CLIP score on the MJHQ30K dataset.

Method	Landscape			People			Food		
	Latency	FID↓	CLIP↑	Latency	FID↓	CLIP↑	Latency	FID↓	CLIP↑
Infinity-2B	1.67	44.1	0.267	1.71	58.91	0.281	1.69	84.2	0.270
+FastVAR	0.60	45.1	0.264	0.61	71.8	0.274	0.61	84.7	0.273
+SkipVAR	1.01	58.1	0.260	1.06	73.7	0.253	0.88	102.3	0.256
+ToProVAR	<b>0.50</b>	<b>44.5</b>	<b>0.264</b>	<b>0.48</b>	<b>58.84</b>	<b>0.283</b>	<b>0.46</b>	<b>84.3</b>	<b>0.274</b>

### Ablation study of the Three-Dimensional Progressive Manipulation Framework on Infinity-2B.

Method	Latency(s)↓	Speed ↑	GenEval↑
Infinity-2B	2.10	-	0.690
+ Scale Depth Loc.	0.47	4.5 ×	0.477
++ Layer Repr. Ident.	0.57	3.7 ×	0.679
+++ Fine-grained Token Prun.	0.61	3.4 ×	0.690

## Visualization



## Reference

- [1] Zhuokun Chen, Jugang Fan, Zhuowei Yu, Bohan Zhuang, and Minghui Tan. Frequency-aware autoregressive modeling for efficient high-resolution image synthesis. 2025a.
- [2] 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.
- [3] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis. arXiv preprint arXiv:2412.04431, 2024.
- [4] 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, 2025a.
- [5] Ziran Qin, Youru Lv, Mingbao Lin, Zeren Zhang, Danping Zou, and Weiyao Lin. Head-aware kv cache compression for efficient visual autoregressive modeling. 2025.