

## Simple Table Test

Type	Model	FID↓	IS↑	Time
GAN	BigGAN [13]	8.43	177.9	—
Diff.	ADM [26]	23.24	101.0	—
Diff.	DiT-XL/2 [63]	3.04	240.8	81
Mask.	MaskGIT [17]	7.32	156.0	0.5†
AR	VQGAN [30]	26.52	66.8	25†
VAR	VAR-d36	2.63	303.2	1

This table shows performance comparison of different generative models.

## Massive Modeling: Scalable Image via Next-Scale Prediction

Zehuan Yuan<sup>2,\*</sup>, Bingyue Peng<sup>2</sup>, Liwei Wang<sup>1,\*</sup>

<sup>1</sup>Peking University <sup>2</sup>Bytedance Inc

keyutian@stu.pku.edu.cn, jiangyi.enjoy@bytedance.com,  
yuanzehuan@bytedance.com, bingyue.peng@bytedance.com, wanglw@pku.edu.cn

Try and explore our online demo at: <https://var.vision>

Codes and models: <https://github.com/FoundationVision/VAR>



Figure 1: **Generated samples from Visual AutoRegressive (VAR) transformers trained on ImageNet.** We show  $512 \times 512$  samples (top),  $256 \times 256$  samples (middle), and zero-shot image editing results (bottom).

### Abstract

We present Visual AutoRegressive modeling (VAR), a new generation paradigm that redefines the autoregressive learning on images as coarse-to-fine “next-scale prediction” or “next-resolution prediction”, diverging from the standard raster-scan “next-token prediction”. This simple, intuitive methodology allows autoregressive (AR) transformers to learn visual distributions fast and can generalize well: VAR, for the *first time*, makes GPT-style AR models surpass diffusion transformers in image generation. On ImageNet  $256 \times 256$  benchmark, VAR significantly improve AR baseline by improving Fréchet inception distance (FID) from 18.65 to 1.73, inception score (IS) from 80.4 to 350.2, with  $20\times$  faster inference speed. It is also empirically verified that VAR outperforms the Diffusion Transformer (DiT) in multiple dimensions including image quality, inference speed, data efficiency, and scalability. Scaling up VAR models exhibits clear power-law scaling laws similar to those observed in LLMs, with linear correlation coefficients near  $-0.998$  as solid evidence. VAR further showcases zero-shot generalization ability in downstream tasks including image in-painting, out-painting, and editing. These results suggest VAR has initially emulated the two important properties of LLMs: **Scaling Laws** and **zero-shot** generalization. We have released all models and codes to promote the exploration of AR/VAR models for visual generation and unified learning.

\*Corresponding authors: wanglw@pku.edu.cn, yuanzehuan@bytedance.com; †: project lead