# Visual Generation Tuning

Jiahao Guo[1,4*‡], Sinan Du[2,4*‡], Jingfeng Yao[1], Wenyu Liu[1], Bo Li[4], Haoxiang Cao[3,4‡]

Kun Gai[4], Chun Yuan[2], Kai Wu[4†], Xinggang Wang[1✉]

[1]Huazhong University of Science and Technology (HUST), [2]Tsinghua University
[3]School of Artificial Intelligence, South China Normal University, [4]Kolors Team, Kuaishou Technology
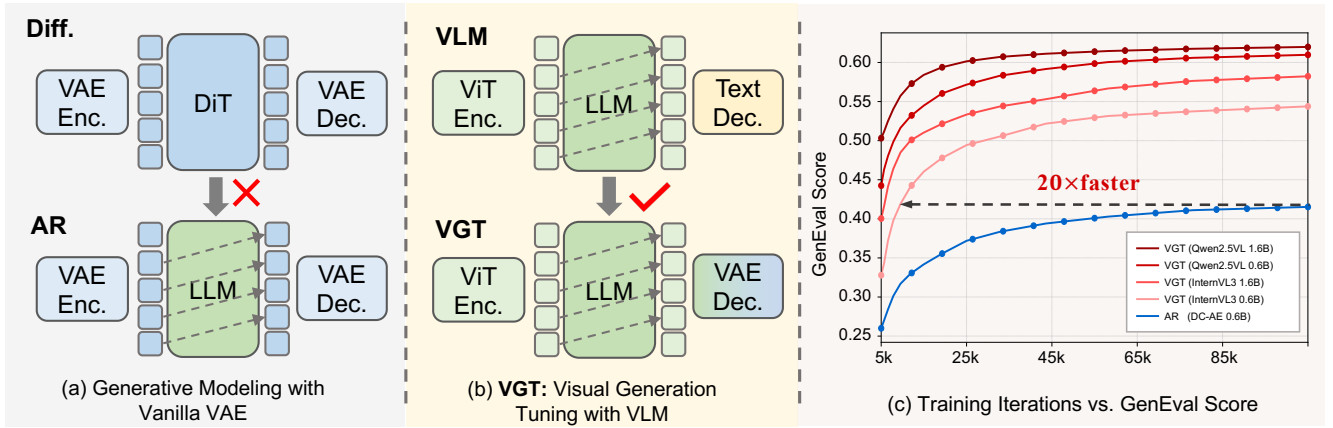
Figure 1. Comparisons of autoregressive generation architectures. (a) Prior continuous autoregressive models directly regress VAE latents produced by pixel-level diffusion VAEs. (b) Our VGT pre-aligns the pretrained VLM semantic encoder with the latent space of the pixel decoder, eliciting emergent generation capabilities. (c) VGT achieves up to **20×** faster training than the VAE-based autoregressive generation (DC-AE) baseline.

## Abstract

*Large Vision Language Models (VLMs) effectively bridge the modality gap through extensive pretraining, acquiring sophisticated visual representations aligned with language. However, it remains underexplored whether these representations, optimized for multimodal understanding tasks, harbor an inherent potential for visual generation. In this paper, we propose **VGT**, Visual **G**eneration **T**uning, a novel paradigm designed to stimulate the underlying capabilities of visual generation within **any** vision language models. By performing efficient visual generation tuning on well-pretrained VLMs, we significantly mitigate the alignment costs and accelerate the convergence of autoregressive modeling in the continuous space(**20×** speedup). Specifically, we dismiss the entangled pixel-level VAEs designed for diffusion transformers and formulate **VGT-AE** through aligning the semantic encoders from pretrained VLMs with*

---

*the latent representations of pixel decoders. In image reconstruction tasks, we achieve **26.67 PSNR** and **0.50 rFID** at a **28×** compression ratio, outperforming specialized VAEs; In visual generation tasks, we achieve state-of-the-art outcomes among autoregressive models, **0.77** on GenEval and **78.73** on DPG-Bench. Furthermore, our proposed VGT showcases significant scaling promise and is versatile for endowing any VLMs trained for multimodal understanding with the capabilities of visual generation, which paves the new avenue to explore next-generation unified multimodal foundation models.Models and codes are available at https://github.com/hustvl/VGT.*

## 1. Introduction

Autoregressive modeling [1, 2, 20, 32, 46, 68] has emerged as a dominant paradigm in language and multimodal generation, demonstrating promising potential for application in visual generation. Previous autoregressive visual generation models primarily relied on vector quantization [11, 34, 53] codebooks to convert images into discrete tokens, inher-

---

* Equal Contibution. † Project Lead. ✉Corresponding Authors.
‡ Work done during internship in Kolors Team, Kuaishou Technology.

ently introducing quantization errors. MAR [26] addresses this through a lightweight flow matching head designed to predict the continuous latent representations of VAEs [39].

However, existing methods [13, 21, 26, 50] neglect a fundamental dilemma: the latent representations learned by vanilla VAEs are poorly aligned with autoregressive modeling. Since VAEs are optimized for pixel-level reconstruction, their latent representations lack semantic structure, which leads to training instability and frequent variance collapse [21, 47, 50]. Inspired by recent studies [42, 48, 67] highlighting that structured semantic representations (*e.g.*, CLIP [37] and DINOv2 [3]) can improve the stability and efficiency in training diffusion transformers (DiTs), we explore a critical question: Vision language models (VLMs) are inherently well-aligned with discriminative visual representations; ***Could we directly leverage this property to transfer visual generation capabilities into well-pretrained VLMs, thereby achieving unified models capable of both multimodal understanding and visual generation?***

In this paper, we propose **VGT**, **V**isual **G**eneration **T**uning, a novel paradigm designed to elicit emergent capabilities of visual generation inherently in ***any*** VLMs originally trained for multimodal understanding tasks. VGT substantially lowers the costs associated with alignment while simultaneously promoting faster convergence in continuous-space autoregressive modeling. Specifically, we formulate **VGT-AE**, which builds upon the semantic encoder of well-pretrained VLMs and trains a pixel decoder for fine-grained reconstruction using a two-stage training strategy. In the first stage, the semantic structure of latent features is preserved via self-distillation loss. Subsequently, channel normalization and noise regularization [47, 50] are applied to enhance the robustness of the latent representations and their adaptability for generative tuning. Inspired by order-agnostic autoregressive methods [13, 26, 33, 63] that facilitate global modeling, we introduce a position-query mechanism that maintains the autoregressive formulation during training while enabling partially parallel decoding during inference. VGT is highly versatile and compatible with various VLMs, *e.g.*, InternVL3 [68]. Moreover, as illustrated in Figure 1c, our VGT demonstrates significant data efficiency compared to vanilla VAE-based autoregressive models. Our contributions are summarized as follows:

1. We propose a novel visual generation tuning (**VGT**) paradigm that enables the generation ability inherent in any VLM originally trained on multimodal understanding tasks.
2. We propose **VGT-AE**, which compresses the semantic features from the VLM vision encoder into a compact structured latent representation for high-fidelity reconstruction and continuous AR modeling.
3. Extensive experiments demonstrate the versatility of our

VGT paradigm, and we achieve SOTA performance in reconstruction (**26.67 PSNR and 0.50 rFID**) and generation (**0.77 GenEval and 78.73 DPG-Bench**) with significant data efficiency (**20x speedup**).

## 2. Related Work

### 2.1. Visual Tokenizers for Generative Modeling

**Discrete.** Early autoregressive visual generation methods predominantly employed vector quantization (VQ) techniques to convert images into discrete token sequences. VQ-VAE [53] pioneered this approach, with subsequent improvements including VQGAN [11], which incorporated adversarial and perceptual losses. Follow-up works enhanced codebook utilization and semantic alignment through hierarchical quantization [69], multi-codebook schemes [70], and bit-efficient representations [41]. Methods such as VQKD [34], TokenFlow [36], and UniTok [30] further integrated semantic understanding into tokenizer training, though they face inherent trade-offs between reconstruction fidelity and semantic comprehension.

**Continuous.** In diffusion models, variational autoencoders (VAEs) [23] serve as the primary framework for learning continuous latent representations. Enhanced variants including KL-16 [39], FluxVAE [24], and VA-VAE [60, 61] optimize the balance between compression efficiency and reconstruction quality. Recent approaches have explored leveraging pretrained visual representations: REPA [64] and REPA-E [25] use semantic features from semantic encoder to guide the denoising process and improve VAE representations, while RAE [67], ClipGen [56], and SVG [42] directly employ semantically rich pretrained features for generation. However, these methods predominantly operate within diffusion frameworks using high-dimensional representations, making them unsuitable for efficient integration with autoregressive modeling. We bridge this gap by leveraging semantic encoders from well-pretrained VLMs to align with low-dimensional latent representations of VAEs for continuous autoregressive generation.

### 2.2. Autoregressive visual generation

**Next-Scale Prediction** methods such as VAR [51] employ a coarse-to-fine generation strategy across multiple resolution scales, utilizing bidirectional context modeling within each scale to capture both local details and global structure. **Next-Set Prediction** methods, commonly implemented as masked generative modeling, include MaskGIT [4] and MAR [26]. These approaches predict subsets of tokens in parallel using bidirectional attention mechanisms, achieving favorable trade-offs between generation efficiency and sample quality through iterative refinement processes. **Next-token Prediction** methods follow the causal autoregressive paradigm of language models,
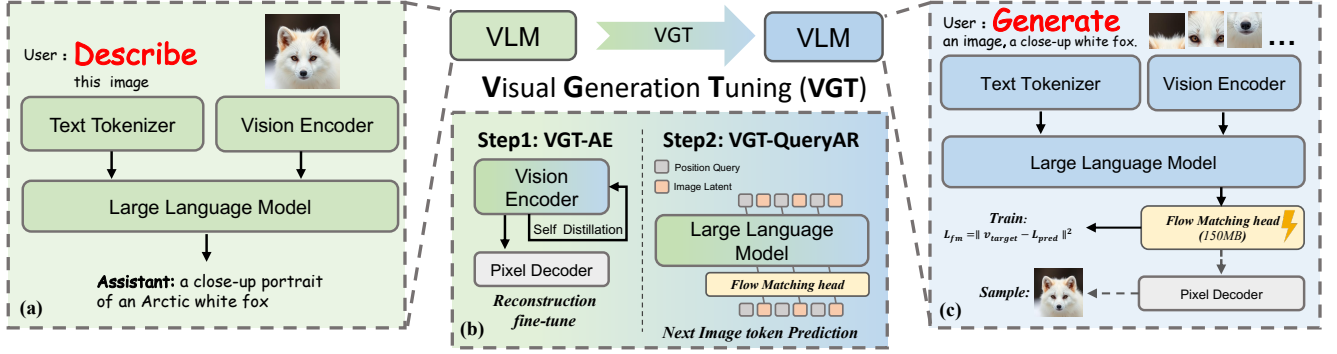
Figure 2. **An overview of our VGT paradigm, which transforms a pre-trained Vision-Language Model (VLM) into a powerful generative model.** (a) Pre-trained VLMs excel at visual understanding by aligning semantic vision encoders with Language Models (LLMs), hinting at their latent image generation capabilities. (b) Our VGT-AE (Visual Generation Tuning-AutoEncoder) aligns the VLM's semantic encoder through reconstruction training. (c) We perform visual generation tuning(VGT) on the VLM by predicting these semantically aligned image latents using a lightweight flow matching head, enabling efficient continuous autoregressive image generation.

generating images token by token in sequential order. Early approaches like VQGAN [11] and LlamaGen [46], used discrete tokens with raster-scan ordering, while recent methods including Ming-UniVision [19], and NextStep-1 [50] have explored continuous token representations. Within this paradigm, researchers have investigated various design choices including discrete versus continuous token representations [47, 52] and fixed raster order versus randomized generation orders [33, 63]. However, these continuous autoregressive methods predominantly borrow the VAEs designed for pixel-level diffusion models, negelecting the misalignment with autoregressive modeling.

## 3. Method

### 3.1. Motivation

We further elucidate the key motivation behind the design of our visual generation tuning methods through pilot experiments. Previous studies on continuous autoregressive modeling typically adopt the pretrained VAEs derived from pixel diffusion models [13, 21, 26, 50]. However, they neglect the representation misalignment with autoregressive models. Inspired by the acceleration of divergence caused by semantic alignment with the VAE latents in diffusion models [42, 64, 67], we perform experiments to validate the conclusion on pretrained VLMs. As presented in Tab. 1 Exp. 1-2, we started by fine-tuning the semantic encoders of pretraiend VLMs without feature constraints, which showcase superior capabilities of reconstruction with significant degradation of semantic structures as evidenced on the multimodal understanding benchmarks. However, a simple self-distillation loss constrains the distribution shift towards fine-grained reconstruction, while improving the understanding and generation. It indicates that the semantic structured representations facilitate the convergence of continuous AR modeling. Moreover, as illustrated in Fig-
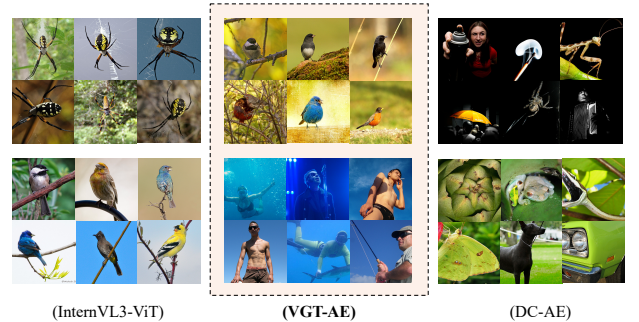


Figure 3. **Visualization of clustering results with different visual representations.** Our VGT-AE is capable of preserving semantic structure while retaining fine-grained textures.

ure 3, we analyze the cluster results of representations of different encoders, which show that the semantic encoder (*e.g., InternViT*) tends to focus on categories (high-level) while pixel encoder (*e.g., DC-AE*) clusters similar textures (low-level). Our VGT-AE achieves an effective compromise, producing representations that are both semantically structured and texturally rich, leading to coherent generalized clusters with refined visual detail. These observations motivate us to enhance the capabilities of visual generation by aligning the semantic encoder with the latent space of the pixel decoder for autoregressive modeling.

### 3.2. VGT-AE

The performance of continuous autoregressive models critically depends on the quality of visual tokenizers. Existing methods face a fundamental trade-off: traditional VAEs lack semantic supervision and suffer from variance collapse [21, 47, 50], while directly employing structured semantic features compromises reconstruction fidelity and yields excessively high-dimensional representations [48, 67], making them difficult to predict with lightweight flow

| # Exp. | Self-Distillation | Reconstruction | | Understanding | | | | Generation | |
|---|---|---|---|---|---|---|---|---|---|
| | | rFID↓ | PSNR↑ | MME-P↑ | MMB↑ | AI2D↑ | TQA↑ | GenEval↑ | DPG-Bench↑ |
| 1 | ✗ | **0.66** | **26.26** | 550.8 | 24.8 | 45.4 | 4.4 | 0.43 | 65.67 |
| 2 | ✓ | 1.13 | 23.95 | **1495.6** | **72.6** | **69.6** | **73.9** | **0.49** | **66.90** |

Table 1. **Pilot experiments that indicate the potential trade-off between understanding, reconstruction and generation.** We observe that over-optimizing the autoencoder towards reconstruction leads to a degradation in generative capability, while employing a self-distillation loss to constrain the latent representations to preserve semantic structures can enhance robustness in generative modeling. Previous works validate this insight under the settings of DiTs while we further demonstrate it in the continuous AR paradigm.

matching heads. To address these limitations, we propose VGT-AE, a semantically aligned visual tokenizer trained through a two-stage paradigm that progressively optimizes for both reconstruction quality and autoregressive compatibility.

**Architecture.** VGT-AE consists of three components: (1) a semantic encoder $\mathcal{E}_{\text{vlm}}$ from a pretrained VLM (*e.g.*, InternViT(vision encoder of InternVL3) [68], QwenViT(vision encoder of Qwen2.5-VL) [2]), (2) a residual projection module $\phi$ that compresses features from dimension $d$ to $d_z = 32$, and (3) a pixel decoder $\mathcal{D}$ adapted from DC-AE [7]. The complete forward pass is:

$$\mathbf{f} = \mathcal{E}_{\text{vlm}}(\mathbf{x}), \quad \mathbf{z} = \phi(\mathbf{f}), \quad \hat{\mathbf{x}} = \mathcal{D}(\mathbf{z}). \quad (1)$$

**Stage 1: Semantic-Preserving Reconstruction.** In the first stage, we jointly optimize the encoder and decoder using a composite loss that combines pixel reconstruction with semantic self-distillation:

$$\mathcal{L}_{\text{rec}} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \mathcal{L}_{\text{LPIPS}}(\mathbf{x}, \hat{\mathbf{x}}) + \mathcal{L}_{\text{GAN}}(\mathbf{x}, \hat{\mathbf{x}}), \quad (2)$$

$$\mathcal{L}_{\text{distill}} = \|\mathcal{E}_{\text{teacher}}(\mathbf{x}) - \mathcal{E}_{\text{vlm}}(\mathbf{x})\|_2^2, \quad (3)$$

$$\mathcal{L}_{\text{stage1}} = \mathcal{L}_{\text{rec}} + \lambda_{\text{distill}} \cdot \mathcal{L}_{\text{distill}}, \quad (4)$$

where $\|\cdot\|_2^2$ denotes mean squared error, $\mathcal{L}_{\text{LPIPS}}$ [65] is the perceptual loss, $\mathcal{L}_{\text{GAN}}$ is the adversarial loss, $\mathcal{E}_{\text{teacher}}$ is the frozen pretrained encoder serving as the distillation target, and $\lambda_{\text{distill}}$ controls distillation strength(set to 1.0 in our experiments). This stage achieves high-fidelity reconstruction while distilling semantic structure into the compact latent space.

**Stage 2: Latent Space Regularization.** The first-stage latent distribution, while semantically structured, does not conform to the standard Gaussian prior (zero mean, unit variance) required for effective flow-based generation. Prior work [21, 47, 50] demonstrates that unregularized latent spaces pose significant challenges for flow-based methods, as mapping from simple noise distributions to complex, unconstrained target distributions becomes prohibitively difficult. Our ablation studies in Table 5 confirm this limitation.

To address this, we freeze the encoder $\mathcal{E}_{\text{vlm}}$ while optimizing only the decoder $\mathcal{D}$ and projection module $\phi$, applying channel-wise layer normalization with Gaussian noise injection:

$$\mathbf{z}_{\text{norm}} = \frac{\mathbf{z} - \mu}{\sigma}, \quad \mathbf{z}_{\text{noisy}} = \mathbf{z}_{\text{norm}} + \epsilon, \quad (5)$$

where $\mu$ and $\sigma$ represent channel-wise mean and standard deviation computed across the batch, and $\epsilon \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$ denotes injected Gaussian noise with $\sigma_{\text{noise}} = 0.1$ in our experiments. The training objective simplifies to:

$$\mathcal{L}_{\text{stage2}} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \mathcal{L}_{\text{LPIPS}}(\mathbf{x}, \hat{\mathbf{x}}). \quad (6)$$

This regularization ensures the resulting latents maintain semantic coherence while becoming more amenable to autoregressive learning, yielding VGT-AE as a semantically meaningful and distributionally stable visual tokenizer.

### 3.3. Autoregressive Modeling

Building upon the semantically structured latents from VGT-AE, we introduce QueryAR for autoregressive visual generation. While conventional raster-scan ordering is prone to error accumulation [47] and existing random-order techniques [26] enhance robustness at the cost of non-autoregressive mask modeling, QueryAR achieves flexible generation order within a standard autoregressive framework through explicit position queries.

**Causal Modeling with Position Queries.** Given a latent sequence $\mathbf{Z} = \{\mathbf{z}_1, \ldots, \mathbf{z}_N\}$ and a random permutation $\pi$, we construct the input sequence by interleaving position queries with corresponding latents: $[Q_{\pi(1)}, \mathbf{z}_{\pi(1)}, Q_{\pi(2)}, \mathbf{z}_{\pi(2)}, \ldots]$, where $Q_i \in \mathbb{R}^{d_{\text{model}}}$ are learnable position embeddings. The language model learns the conditional distribution:

$$p_\theta(\mathbf{z}_{\pi(t)} \mid \mathbf{H}_{<t}, Q_{\pi(t)}), \quad (7)$$

where $\mathbf{H}_{<t}$ denotes the causal context containing all previously generated latents and their associated position queries. This formulation maintains autoregressive properties while supporting flexible generation orders.

**Continuous Latent Modeling via Flow Matching.** For modeling continuous latents, we employ a lightweight flow matching head [26] conditioned on language model hidden states. The language model processes input sequences to produce hidden states $\mathbf{H} = f_\theta([Q_{\pi(1)}, \mathbf{z}_{\pi(1)}, \ldots])$, which

serve as conditioning signals for the flow matching procedure. The training objective follows the standard formulation:

$$\mathcal{L}_{\text{fm}} = \mathbb{E}_{t,\epsilon} \left[ \|(\mathbf{z}_{\text{target}} - \epsilon) - v_\theta(\mathbf{z}_t, t, \mathbf{H})\|_2^2 \right], \quad (8)$$

where $\mathbf{z}_t = t \cdot \mathbf{z}_{\text{target}} + (1 - t) \cdot \epsilon$ is the linearly interpolated latent at timestep $t \in [0, 1]$, $\epsilon \sim \mathcal{N}(0, I)$ is Gaussian noise, and $v_\theta$ is the predicted vector field. For timestep scheduling, we implement the dimensional-adaptive strategy from [12] with $t_m = \frac{\alpha t_n}{1 + (\alpha - 1)t_n}$, where $\alpha = \sqrt{m/n}$ and $n = 4096$ serves as the reference dimension.

**Parallel Decoding for Efficient Inference.** The position-query mechanism substantially enhances inference efficiency through partially parallel decoding. Given previously generated latents $\mathbf{z}_{1:k}$ and target position queries $Q_{k+1:k+m}$, the language model processes these inputs to produce hidden states for the subsequent $m$ positions in a single forward pass:

$$\mathbf{H}_{k+1:k+m} = f_\theta(\mathbf{z}_{1:k}, Q_{k+1:k+m}). \quad (9)$$

These hidden states then condition the flow matching head for deterministic sampling:

$$\hat{\mathbf{Z}}_{k+1:k+m} = \text{FlowSample}(\mathbf{H}_{k+1:k+m}, \epsilon_\theta), \quad (10)$$

where $\epsilon_\theta$ denotes the flow matching head. Our ablation study in Section 4.4.5 validates this design, demonstrating that multiple tokens can be generated simultaneously while preserving generation quality through the position-query mechanism.

# 4. Experiments

## 4.1. Experimental Setups

**Datasets.** VGT is pretrained on the BLIP3-o [8] open-source dataset, which comprises 27M samples recaptioned by Qwen2.5-VL-7B [2], 5M samples from CC12M [5], and 4M synthesized images from JourneyDB [45]. For supervised fine-tuning, we use a total of 200K high-quality aesthetic samples collected from ShareGPT-4o [10], Echo-4o [62], and BLIP3-o [8], to enhance visual quality and prompt alignment.

**Implementation Details.** We develop VGT models with 0.6B and 1.6B parameters, instantiated from the Qwen2.5-VL [2] and InternVL3 [68] VLM families. VGT-AE adopts the vision encoders of Qwen2.5-VL (QwenViT) and InternVL3 (InternViT) and pairs them with the DC-AE decoder [7]. For autoregressive visual generation, QueryAR uses VGT-AE as the visual tokenizer and shares the language model of the underlying VLM. To achieve target model sizes, we apply pruning to the Qwen2.5-VL models. We pretrain using the AdamW optimizer with a batch size of

| Method | Ratio | rFID↓ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|
| *Generative Only Tokenizer* | | | | |
| VQGAN [11] | 16 | 4.98 | 20.00 | 0.629 |
| LlamaGen [46] | 16 | 2.19 | 20.79 | 0.675 |
| VAR [51] | 16 | 1.00 | 22.63 | 0.755 |
| Open-MAGVIT2 [29] | 16 | 1.67 | 22.70 | 0.640 |
| RAE [67] | 16 | 0.49 | 19.23 | 0.620 |
| SD-VAE [39] | 16 | 2.64 | 22.13 | 0.590 |
| DC-AE [7] | 32 | 0.69 | 23.85 | 0.660 |
| *CLIP-based Tokenizer* | | | | |
| VILA-U [58] | 16 | 1.80 | - | - |
| TokenFlow [36] | 16 | 1.37 | 21.41 | 0.687 |
| DualViTok [44] | 16 | 1.37 | 22.53 | 0.741 |
| UniLIP [48] | 32 | 0.79 | 22.99 | 0.747 |
| *Ours* | | | | |
| **VGT-AE (QwenViT)** | 28 | 1.93 | 20.12 | 0.677 |
| **VGT-AE (InternViT)** | 28 | **0.50** | **26.67** | **0.863** |

Table 2. Comparisons of reconstruction quality on the $256 \times 256$ ImageNet 50k validation set. Ratio indicates the downsample ratio relative to the original image. Best results are highlighted in bold. VGT-AE (QwenViT) and VGT-AE (InternViT) use the vision encoders of Qwen2.5-VL and InternVL3, respectively.

256 and a learning rate of $2.0 \times 10^{-4}$, with cosine learning-rate decay. The optimizer settings are $\beta = (0.9, 0.95)$ and weight decay 0.05. EMA is applied with a decay rate of 0.9999. Models are pretrained for 100K steps, and ablation experiments are conducted for 50K steps. Fine-tuning for autoregressive visual generation is performed for 5,000 iterations with a learning rate of $5 \times 10^{-5}$. We use a causal mask with a maximum sequence length of 1024 tokens.

**Evaluation Metrics.** We assess reconstruction quality using rFID, PSNR, and SSIM on the ImageNet-1K validation set. For visual generation, we evaluate on GenEval [17] and DPG-Bench [18]. For multimodal understanding, we evaluate on MME-P: MME-Perception [14], MMB: MMBench-en [28], TQA: TextVQA [43] and AI2D [22].

## 4.2. Visual Tokenizer

VGT-AE, a semantically aligned visual tokenizer, successfully bridges the powerful representational capabilities of vision encoders from VLMs with the demands of high-fidelity image reconstruction through a carefully designed two-stage training strategy. Experimental results in Tab. 2 demonstrate that our two-stage approach enables high-quality reconstruction, achieving state-of-the-art reconstruction performance. Notably, VGT-AE-InternViT attains remarkable metrics on the ImageNet $256\times256$ validation set, with rFID of **0.50**, PSNR of **26.67**, and SSIM of **0.863**, significantly surpassing existing generative-only and CLIP-based tokenizers. We provide qualitative visualizations in Fig. 4.

| Method | Data. | # Params. | GenEval [17] | | | | | DPG-Bench [18] | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Two Object | Counting | Colors | Position | Overall↑ | Global | Attribute | Relation | Other | Overall↑ |
| *Diffusion-based Model* | | | | | | | | | | | | |
| SDv1.5 [40] | >2000M | 0.9B | 0.38 | 0.35 | 0.76 | 0.04 | 0.43 | 74.63 | 75.39 | 73.49 | 67.81 | 63.18 |
| PixArt-α [6] | 25M | 0.6B | 0.50 | 0.44 | 0.80 | 0.08 | 0.48 | 74.97 | 78.60 | 82.57 | 76.96 | 71.11 |
| SDv2.1 [40] | - | 0.9B | 0.51 | 0.44 | 0.85 | 0.07 | 0.50 | - | - | - | - | - |
| SDXL [35] | - | 2.6B | 0.74 | 0.39 | 0.85 | 0.15 | 0.55 | 83.27 | 80.91 | 86.76 | 80.41 | 74.65 |
| DALLE3 [38] | - | - | 0.87 | 0.47 | 0.83 | **0.43** | 0.67 | **90.97** | 88.39 | **90.58** | **89.83** | 83.50 |
| SD3-Medium [40] | - | 2B | **0.94** | **0.72** | **0.89** | 0.33 | **0.74** | 87.90 | **88.83** | 80.70 | 88.68 | **84.08** |
| *Autoregressive-based Model* | | | | | | | | | | | | |
| Chameleon [49] | >1.4B | 7B | - | - | - | - | 0.39 | - | - | - | - | - |
| Fuild [13] | 2048M | 0.4B | - | - | - | - | 0.45 | - | - | - | - | - |
| Fuild [13] | 2048M | 0.7B | - | - | - | - | 0.51 | - | - | - | - | - |
| LlamaGen [46] | - | 0.8B | 0.34 | 0.21 | 0.58 | 0.07 | 0.32 | 81.76 | 76.17 | 84.76 | 58.40 | 64.84 |
| EMU3-Gen [55] | - | 8B | 0.71 | 0.34 | 0.81 | 0.17 | 0.54 | 85.21 | 86.84 | **90.22** | 83.15 | 80.60 |
| TokenFlow [36] | - | 13B | 0.66 | 0.40 | 0.84 | 0.17 | 0.55 | 78.72 | 81.29 | 85.22 | 71.20 | 73.38 |
| SEED-X [15] | - | 0.7B | 0.65 | 0.31 | 0.80 | 0.18 | 0.51 | - | - | - | - | - |
| Janus [57] | 198M | 1.3B | 0.68 | 0.30 | 0.84 | 0.46 | 0.61 | 82.33 | 87.70 | 85.46 | 86.41 | 79.68 |
| SimpleAR [54] | - | 1.5B | **0.90** | - | - | 0.28 | 0.63 | **87.97** | - | 88.33 | - | 81.97 |
| VAR [51] | - | - | - | - | - | - | 0.53 | - | - | - | - | 71.08 |
| Janus-Pro [9] | 198M | 1B | 0.82 | 0.51 | **0.89** | 0.65 | 0.73 | 87.58 | **88.17** | 88.98 | **88.30** | 82.63 |
| NextStep-1 [50] | >2048M | 14B | - | - | - | - | 0.63 | - | - | - | - | **85.28** |
| **VGT (InternVL3)** | <25M | 0.6B | 0.76 | 0.51 | 0.78 | 0.70 | 0.71 | 81.55 | 82.96 | 85.91 | 82.21 | 72.01 |
| **VGT (InternVL3)** | <25M | 1.6B | 0.83 | **0.61** | 0.84 | 0.70 | 0.75 | 86.55 | 83.76 | 84.36 | 85.14 | 74.43 |
| **VGT (Qwen2.5-VL)** | <25M | 0.6B | 0.78 | 0.54 | 0.83 | 0.68 | 0.72 | 85.16 | 84.89 | 85.50 | 82.00 | 75.18 |
| **VGT (Qwen2.5-VL)** | <25M | 1.6B | 0.85 | 0.59 | 0.87 | **0.74** | **0.77** | 87.41 | 85.67 | 86.78 | 87.73 | 78.73 |

Table 3. Comparisons of visual generation quality on GenEval [17] and DPG-Bench [18]. Bold numbers indicate the best performance in each column. Our VGT models demonstrate highly competitive results, with the VGT-Qwen2.5-VL variant establishing new state-of-the-art performance on GenEval while trained on remarkably limited data (only 25M samples).
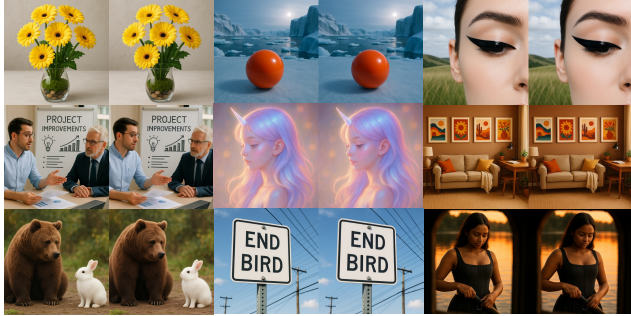


Figure 4. Visualization of reconstruction results from our VGT-AE-InternViT. **Left**: input image; **Right**: reconstructed image.

## 4.3. Visual Generation

We systematically evaluate the visual generation performance of our proposed VGT paradigm on two established benchmarks: GenEval [17] and DPG-Bench [18]. As shown in Table 3, our method demonstrates strong performance, particularly on the GenEval benchmark. As demonstrated in Figure 5, our VGT paradigm effectively endows well-pretrained vision–language models with the ability to produce diverse and realistic high-quality images. The VGT model based on Qwen2.5-VL [2] achieves a state-of-the-art overall score of **0.77**, surpassing existing autoregressive generation models including Janus-Pro [9], TokenFlow [36], and SimpleAR [54]. Notably, it also competes favorably with or even exceeds several large diffusion models such as SDXL [35], SD3-Medium [40], and

DALLE3 [38] across multiple sub-metrics. This strong performance is achieved with only 25M training samples, highlighting VGT's superior instruction-following ability and its effectiveness in constrained generation. A key factor behind this improvement is our visual representation module, VGT-AE, which builds on representations learned by large vision–language models. This design offers inherent vision–language alignment, allowing the generator to maintain strong semantic consistency with textual instructions while avoiding the structural information loss commonly seen in traditional VAEs.
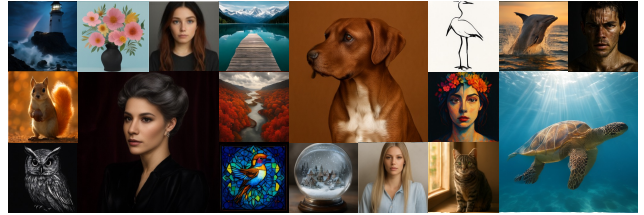


Figure 5. Our **VGT-1.6B** based on Qwen2.5-VL and InternVL3, endow pretrained vision–language models trained on multimodal understanding tasks with high-quality visual generation, enabling them to produce diverse and realistic images.

Our results challenge long-standing assumptions in visual generation. Conventional autoregressive models typically require orders of magnitude more training data (hundreds of millions to billions of samples) yet still exhibit unstable performance. In contrast, VGT shows clear advantages on challenging constraint-based tasks such as *Count-*

*ing* and *Position*. More importantly, VGT achieves these improvements while maintaining substantially higher training efficiency compared to diffusion-based approaches. For instance, SD3-Medium [40] relies on 2B parameters and expensive training procedures to achieve high GenEval scores, whereas our VGT-Qwen2.5-VL-1.6B model attains comparable performance using only 25M training samples. This challenges the long-held belief that autoregressive models inherently produce lower-quality images than diffusion models of comparable scale. Even our smaller VGT-0.6B model surpasses larger systems such as NextStep-1 [50] and Janus-Pro [9] on several metrics, further confirming that the vision–language aligned representation in VGT enables stronger generalization and more controllable visual generation. These results point to a promising new paradigm for efficient yet high-quality autoregressive visual generation.

## 4.4. Ablation Studies

To comprehensively evaluate the design choices in VGT, we conduct ablation studies across five critical components: the VGT-AE decoder architecture (Section 4.4.1), training methodology (Section 4.4.2), the reconstruction-generation trade-off (Section 4.4.3), cross-architecture compatibility (Section 4.4.4), and QueryAR structure (Section 4.4.5).

### 4.4.1. VGT-AE Architecture

| Decoder Type | rFID↓ | PSNR↑ | SSIM↑ | Param. |
|---|---|---|---|---|
| *VGT-AE with Qwen2.5-VL Vision Encoder* | | | | |
| ViT | 17.32 | 18.56 | 0.518 | 845M |
| SD-VAE | 2.81 | 19.22 | 0.640 | **743M** |
| DC-AE† | 2.71 | 19.41 | **0.647** | 853M |
| **DC-AE** | **2.67** | **19.68** | 0.563 | 853M |
| *VGT-AE with InternVL3 Vision Encoder* | | | | |
| ViT | 15.80 | 20.01 | 0.611 | 464M |
| SD-VAE | 1.32 | 23.10 | **0.816** | **362M** |
| DC-AE† | 1.14 | 23.86 | 0.792 | 472M |
| **DC-AE** | **1.13** | **23.95** | 0.801 | 472M |

Table 4. Ablation study on decoder architectures for VGT-AE with different vision-language model encoders. † indicates models trained without pretrained weights.

In Table 4, we compare several mainstream decoder architectures for VGT-AE: the high-compression DC-AE from SANA [7], the SD-VAE decoder from Stable Diffusion [39], and a ViT-based decoder of comparable parameter scale. Results show that both DC-AE and SD-VAE decoders demonstrate significant advantages regardless of whether Qwen2.5-VL or InternVL3 vision encoders are employed, while the ViT decoder substantially lags behind across all metrics.

For instance, in the InternVL3 configuration, DC-AE reduces rFID from 15.80 (ViT) to **1.13** and improves PSNR to **23.95**. The relatively small performance gap between DC-AE and SD-VAE (1.13 vs. 1.32 rFID) indicates that both

effectively recover image details while maintaining semantic consistency. Based on these findings, we select DC-AE as our default decoder architecture to achieve the optimal balance between semantic structure preservation, visual detail restoration, and parameter efficiency.

### 4.4.2. VGT-AE Training Methodology

| Exp. | Training Strategy | | | rFID↓ | PSNR↑ | Geneval↑ |
|---|---|---|---|---|---|---|
| | Stage | Norm | Noise | | | |
| 1 | Stage1 | ✗ | ✗ | **0.98** | **24.30** | 0.36 |
| 2 | Stage2 | ✓ | ✗ | 1.05 | 24.21 | 0.52 |
| 3 | Stage2 | ✓ | 0.1 | 1.13 | 23.95 | **0.54** |

Table 5. Ablation study on VGT-AE's training strategies. Stage1 focuses on reconstruction, while Stage2 enhances generation through latent space regularization. The performance progression from Experiment 1 to 3 underscores the importance of regularization for learning generation-friendly representations.

Our ablation study provides key findings regarding VGT-AE's training paradigm. As summarized in Table 5, Stage 1 training achieves the best reconstruction quality but yields poor generation performance (GenEval=0.36). This suggests that latents optimized solely for reconstruction are suboptimal for the denoising process in autoregressive generation. Introducing normalization in Stage 2 (Experiment 2) improves generation substantially, and further adding mild noise injection (Experiment 3) yields the best generation result (GenEval=**0.54**), despite a slight degradation in reconstruction metrics.
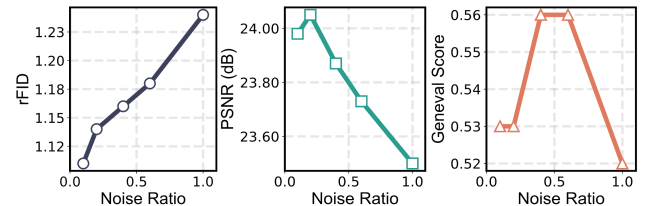


Figure 6. Impact of noise regularization on the reconstruction-generation trade-off. While reconstruction quality (rFID, PSNR) degrades with increasing noise, generation performance (GenEval) peaks at moderate noise intensities ($\sigma = 0.4$–$0.6$).

Figure 6 further elucidates this trade-off. Reconstruction quality decreases monotonically with noise level, whereas generation performance peaks at moderate levels ($\sigma = 0.4$-$0.6$). This demonstrates that an appropriate level of regularization is crucial for striking a balance.

Collectively, these results validate VGT-AE's two-stage training strategy as an effective paradigm for learning latent representations that reconcile high-fidelity reconstruction with effective autoregressive generation.

| Model | rFID↓ | PSNR↑ | Geneval↑ | DPG-Bench↑ |
|---|---|---|---|---|
| VGT(Qwen2.5-VL) | 1.93 | 20.12 | **0.72** | **75.18** |
| VGT(InternVL3) | **0.50** | **26.67** | 0.71 | 72.01 |
| VGT(InternVL3, Low Rec.) | 1.13 | 23.95 | **0.72** | 73.50 |

Table 6. Comparison of reconstruction and generation performance across VGT models(0.6B). Reconstruction-oriented metrics (rFID, PSNR) and generation-oriented benchmarks (Geneval, DPG-Bench) exhibit a consistent inverse trend, highlighting a structural trade-off inherent to autoregressive latent modeling.

### 4.4.3. Reconstruction vs. Generation

We study the trade-off between reconstruction fidelity and generative capability using three VGT-AE variants instantiated with Qwen2.5-VL [2] and InternVL3 [68]. As shown in Table 6, the reconstruction-oriented InternVL3 variant, which yields the most compact and pixel-aligned latent representations, achieves the best reconstruction scores but underperforms on GenEval [17] and DPG-Bench [18]. Relaxing its reconstruction objective (InternVL3 Low Rec.) degrades pixel-level fidelity yet consistently improves generation metrics, and the Qwen2.5-VL-based VGT shows a similar trend favoring generation. Figure 7 further indicates that these gains correlate with more dispersed, semantically separated manifolds, whereas the reconstruction-oriented variant forms dense, entangled clusters. Qualitatively (Figure 8), however, the high-reconstruction InternVL3-based VGT still produces sharper illumination and finer textures than its Qwen2.5-VL counterpart, indicating that higher automatic generation metrics (e.g., GenEval and DPG-Bench) are not necessarily aligned with perceptual sharpness.
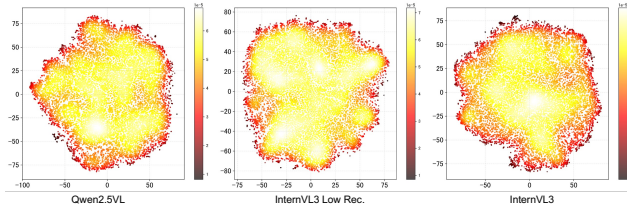


Figure 7. t-SNE of latent tokens from three VGT-AE variants. Generation-oriented models (Qwen2.5-VL and InternVL3 Low Rec.) yield dispersed, semantically structured manifolds, while the reconstruction-oriented InternVL3 forms compact, entangled clusters, reflecting different requirements on the latent space for reconstruction vs. autoregressive generation.

This trade-off echoes findings in diffusion transformers, where semantically structured representations (e.g., CLIP-/DINO-like features) improve stability and sample quality at the expense of exact pixel matching [42, 48, 67], and aligns with observations in discrete visual tokenization [34, 36, 44]. Together, these results highlight that *representation design* is central to visual generation: latents optimized purely for pixel reconstruction are overly constrained for autoregressive prediction, whereas semantically organized, weakly coupled latents are more favorable for generation [50, 61]. VGT-AE targets this regime via a two-stage training scheme: self-distillation to preserve the VLM encoder's semantic structure, followed by channel normalization and noise regularization [47, 50] to enhance robustness and reconstruction quality. This allows VGT to approach the reconstruction fidelity of specialized VAEs while retaining a latent space intrinsically well-suited for continuous-space autoregressive generation.



Figure 8. Qualitative comparison between InternVL3-based VGT (**left**) and Qwen2.5-VL-based VGT (**right**). The InternVL3 variant, despite lower generation metrics, produces noticeably sharper illumination and more detailed textures.

### 4.4.4. Different Language Models in VGT

We study the coupling between the autoencoder (AE) and the language model (LLM) by evaluating four VGT configurations: (i) AE and LLM both from Qwen2.5-VL, (ii) both from InternVL3, (iii) a Qwen2.5-VL AE with an InternVL3 LLM, and (iv) a VAE-based DC-AE baseline. As summarized in Table 7, using AE and LLM from the same VLM family yields the strongest overall performance, reflecting coherent cross-modal alignment from large-scale multimodal pre-training. The mismatched configuration (Qwen2.5-VL AE + InternVL3 LLM) still surpasses the VAE baseline on generation metrics, while the VAE lags behind across all evaluations, confirming that pixel-reconstruction-oriented latents are poorly suited for autoregressive decoding [7, 39].

| Autoencoder | LLM From | Geneval↑ | DPG-Bench↑ |
|---|---|---|---|
| VGT-AE (Qwen2.5-VL) | Qwen2.5-VL | **0.77** | **78.73** |
| VGT-AE (InternVL3) | InternVL3 | 0.75 | 74.43 |
| VGT-AE (Qwen2.5-VL) | InternVL3 (mismatch) | 0.75 | 72.46 |
| DC-AE (VAE baseline) | InternVL3 | 0.64 | 67.20 |

Table 7. Performance comparison across different Autoencoder and LLM combinations within the VGT paradigm. Even in mismatched settings, VGT-AE significantly outperforms the VAE baseline (DC-AE), demonstrating the robustness and generative friendliness of VGT latent representations.

These observations support a central insight of the VGT paradigm: visual representations learned from multimodal understanding data can be transferred to language models for autoregressive visual generation with a *low alignment cost*. Because VGT-AE is trained inside a multimodally aligned VLM (Section 4.2), its latent tokens inherit semantic structure that is naturally compatible with different LLMs, with matched AE–LLM pairs amplifying this effect and mismatched pairs still outperforming VAEs that must be aligned from scratch. Thus, VGT effectively upgrades existing VLMs into unified models for both multimodal understanding and visual generation, and even under imperfect component pairing (Table 7), its semantically organized latent space makes bridging to autoregressive generation substantially easier than with conventional VAE-based approaches.

### 4.4.5. QueryAR vs. MAR

| # Exp. | Method | Accel. Ratio | Geneval↑ | DPGBench↑ |
|--------|--------|--------------|----------|-----------|
| 1 | MAR | 1× | **0.51** | 74.39 |
| 2 | MAR | 4× | 0.46 | **75.29** |
| 3 | QueryAR | 1× | **0.62** | 77.78 |
| 4 | QueryAR | 4× | 0.59 | **78.14** |

Table 8. Comparative analysis of decoding strategies: QueryAR vs. MAR across different acceleration ratios. QueryAR maintains competitive performance even at 4× acceleration, demonstrating its robustness for parallel decoding.

To rigorously evaluate the efficacy of QueryAR, we conduct a comparative analysis against MAR [26], a prominent random-order autoregressive method, using the VGT-AE Qwen2.5-VL model. As shown in Table 8, the performance of MAR degrades under acceleration, with the GenEval score dropping from 0.51 (1×) to 0.46 (4×). This indicates a compromise in overall generation quality under increased acceleration.

In stark contrast, QueryAR consistently outperforms MAR across all tested acceleration ratios. At 1× decoding, QueryAR achieves a superior 0.62 GenEval score. Crucially, at 4× acceleration, QueryAR maintains a strong GenEval of 0.59 and achieves the highest DPG-Bench score (78.14) in this comparison. This robust performance demonstrates QueryAR's ability to achieve significant inference speed-ups without compromising generation quality.

The qualitative results presented in Figure 9 further confirm QueryAR's exceptional robustness and coherence, generating high-quality images even at an aggressive 16× acceleration ratio. This visual evidence corroborates our quantitative findings, confirming that critical visual details and semantic consistency are preserved despite the highly parallelized decoding. Collectively, these results confirm that QueryAR, through its innovative position-query mechanism, successfully learns a more stable and efficient autoregressive latent transition, effectively balancing superior generative performance with enhanced inference efficiency.



Figure 9. **QueryAR** simultaneously generates 4 or 16 tokens while maintaining high-quality image outputs.

## 5. Conclusion

In this work, we unlock the latent generative potential within large vision-language models by introducing VGT, a visual generation tuning paradigm. This approach effectively activates generative capabilities without requiring architectural redesign or incurring prohibitive training costs. Our core innovation, VGT-AE, meticulously aligns pretrained VLM semantic encoders with lightweight pixel decoders. This alignment directly addresses the inherent mismatch between traditional VAEs and autoregressive modeling, facilitating stable and efficient learning within a continuous latent space. Extensive experiments demonstrate that VGT-AE achieves state-of-the-art reconstruction performance even under high compression ratios, and significantly advances autoregressive visual generation. VGT delivers competitive scores on both GenEval and DPG-Bench, all while maintaining fast convergence and strong scalability. Beyond its specific model instance, VGT provides a general blueprint for equipping any multimodal understanding model with high-quality visual generation capabilities. This work paves a promising path toward next-generation unified multimodal foundation models, capable of seamless perception-generation synergy, with future directions potentially extending the VGT framework in line with recent multimodal developments [16, 27, 31, 59, 66].

## 6. Acknowledgements

## References

[1] Lei Bai, Zhongrui Cai, Yuhang Cao, Maosong Cao, Weihan Cao, Chiyu Chen, Haojiong Chen, Kai Chen, Pengcheng Chen, Ying Chen, et al. Intern-s1: A scientific multimodal foundation model. *arXiv preprint arXiv:2508.15763*, 2025. 1

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 4, 5, 6, 8

[3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 2

[4] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11315–11325, 2022. 2

[5] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pretraining to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3558–3568, 2021. 5

[6] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\alpha$: Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2023. 6

[7] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv preprint arXiv:2410.10733*, 2024. 4, 5, 7, 8

[8] Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025. 5

[9] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025. 6, 7

[10] Erfei Cui, Yinan He, Zheng Ma, Zhe Chen, Hao Tian, Weiyun Wang, Kunchang Li, Yi Wang, Wenhai Wang, Xizhou Zhu, Lewei Lu, Tong Lu, Yali Wang, Limin Wang, Yu Qiao, and Jifeng Dai. Sharegpt-4o: Comprehensive multimodal annotations with gpt-4o, 2024. 5

[11] 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. 1, 2, 3, 5

[12] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. 5

[13] Lijie Fan, Tianhong Li, Siyang Qin, Yuanzhen Li, Chen Sun, Michael Rubinstein, Deqing Sun, Kaiming He, and Yonglong Tian. Fluid: Scaling autoregressive text-to-image generative models with continuous tokens. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. 2, 3, 6

[14] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 5

[15] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024. 6

[16] Zigang Geng, Yibing Wang, Yeyao Ma, Chen Li, Yongming Rao, Shuyang Gu, Zhao Zhong, Qinglin Lu, Han Hu, Xiaosong Zhang, et al. X-omni: Reinforcement learning makes discrete autoregressive image generative models great again. *arXiv preprint arXiv:2507.22058*, 2025. 9

[17] 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. 5, 6, 8

[18] 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. 5, 6, 8

[19] Ziyuan Huang, DanDan Zheng, Cheng Zou, Rui Liu, Xiaolong Wang, Kaixiang Ji, Weilong Chai, Jianxin Sun, Libin Wang, Yongjie Lv, et al. Ming-univision: Joint image understanding and generation with a unified continuous tokenizer. *arXiv preprint arXiv:2510.06590*, 2025. 3

[20] Kat Kampf and Nicole Brichtova. Experiment with gemini 2.0 flash native image generation, 2025. *URL https://developers. googleblog. com/en/experiment-with-gemini-20-flash-native-image-generation*, 3, 2025. 1

[21] Guolin Ke and Hui Xue. Hyperspherical latents improve continuous-token autoregressive generation. *arXiv preprint arXiv:2509.24335*, 2025. 2, 3, 4

[22] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images, 2016. 5

[23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2

[24] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. 2

[25] Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng. Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint arXiv:2504.10483*, 2025. 2

[26] 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, 3, 4, 9

[27] Huijie Liu, Bingcan Wang, Jie Hu, Xiaoming Wei, and Guoliang Kang. Omni-dish: Photorealistic and faithful image generation and editing for arbitrary chinese dishes. *arXiv preprint arXiv:2504.09948*, 2025. 9

[28] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024. 5

[29] Zhuoyan Luo, Fengyuan Shi, Yixiao Ge, Yujiu Yang, Limin Wang, and Ying Shan. Open-magvit2: An open-source project toward democratizing auto-regressive visual generation. *arXiv preprint arXiv:2409.04410*, 2024. 5

[30] Chuofan Ma, Yi Jiang, Junfeng Wu, Jihan Yang, Xin Yu, Zehuan Yuan, Bingyue Peng, and Xiaojuan Qi. Unitok: A unified tokenizer for visual generation and understanding. *arXiv preprint arXiv:2502.20321*, 2025. 2

[31] Qingyang Mao, Qi Cai, Yehao Li, Yingwei Pan, Mingyue Cheng, Ting Yao, Qi Liu, and Tao Mei. Visual autoregressive modeling for instruction-guided image editing. *arXiv preprint arXiv:2508.15772*, 2025. 9

[32] AI Open. Gpt-4 technical report, 2023. 1

[33] 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. *arXiv preprint arXiv:2412.01827*, 2024. 2, 3

[34] Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. Beit v2: Masked image modeling with vector-quantized visual tokenizers. *arXiv preprint arXiv:2208.06366*, 2022. 1, 2, 8

[35] 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, 2023. 6

[36] Liao Qu, Huichao Zhang, Yiheng Liu, Xu Wang, Yi Jiang, Yiming Gao, Hu Ye, Daniel K Du, Zehuan Yuan, and Xinglong Wu. Tokenflow: Unified image tokenizer for multimodal understanding and generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2545–2555, 2025. 2, 5, 6, 8

[37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 2

[38] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint*, abs/2204.06125, 2022. 6

[39] 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, 5, 7, 8

[40] 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. 6, 7

[41] Fengyuan Shi, Zhuoyan Luo, Yixiao Ge, Yujiu Yang, Ying Shan, and Limin Wang. Scalable image tokenization with index backpropagation quantization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16037–16046, 2025. 2

[42] Minglei Shi, Haolin Wang, Wenzhao Zheng, Ziyang Yuan, Xiaoshi Wu, Xintao Wang, Pengfei Wan, Jie Zhou, and Jiwen Lu. Latent diffusion model without variational autoencoder. *arXiv preprint arXiv:2510.15301*, 2025. 2, 3, 8

[43] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019. 5

[44] Wei Song, Yuran Wang, Zijia Song, Yadong Li, Haoze Sun, Weipeng Chen, Zenan Zhou, Jianhua Xu, Jiaqi Wang, and Kaicheng Yu. Dualtoken: Towards unifying visual understanding and generation with dual visual vocabularies. *arXiv preprint arXiv:2503.14324*, 2025. 5, 8

[45] Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances in neural information processing systems*, 36:49659–49678, 2023. 5

[46] 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. 1, 3, 5, 6

[47] Yutao Sun, Hangbo Bao, Wenhui Wang, Zhiliang Peng, Li Dong, Shaohan Huang, Jianyong Wang, and Furu Wei. Multimodal latent language modeling with next-token diffusion. *arXiv preprint arXiv:2412.08635*, 2024. 2, 3, 4, 8

[48] Hao Tang, Chenwei Xie, Xiaoyi Bao, Tingyu Weng, Pandeng Li, Yun Zheng, and Liwei Wang. Unilip: Adapting clip for unified multimodal understanding, generation and editing. *arXiv preprint arXiv:2507.23278*, 2025. 2, 3, 5, 8

[49] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024. 6

[50] NextStep Team, Chunrui Han, Guopeng Li, Jingwei Wu, Quan Sun, Yan Cai, Yuang Peng, Zheng Ge, Deyu Zhou, Haomiao Tang, et al. Nextstep-1: Toward autoregressive image generation with continuous tokens at scale. *arXiv preprint arXiv:2508.10711*, 2025. 2, 3, 4, 6, 7, 8

[51] 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, 5, 6

[52] Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. In *European Conference on Computer Vision*, pages 292–309. Springer, 2024. 3

[53] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 1, 2

[54] Junke Wang, Zhi Tian, Xun Wang, Xinyu Zhang, Weilin Huang, Zuxuan Wu, and Yu-Gang Jiang. Simplear: Pushing the frontier of autoregressive visual generation through pre-training, sft, and rl. *arXiv preprint arXiv:2504.11455*, 2025. 6

[55] 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. 6

[56] Zihao Wang, Wei Liu, Qian He, Xinglong Wu, and Zili Yi. Clip-gen: Language-free training of a text-to-image generator with clip. *arXiv preprint arXiv:2203.00386*, 2022. 2

[57] Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified multimodal understanding and generation. *arXiv preprint arXiv:2410.13848*, 2024. 6

[58] Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. Vila-u: a unified foundation model integrating visual understanding and generation. *arXiv preprint arXiv:2409.04429*, 2024. 5

[59] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. *arXiv preprint arXiv:2409.11340*, 2024. 9

[60] Jingfeng Yao, Cheng Wang, Wenyu Liu, and Xinggang Wang. Fasterdit: Towards faster diffusion transformers training without architecture modification. *Advances in Neural Information Processing Systems*, 37:56166–56189, 2024. 2

[61] Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 2, 8

[62] Junyan Ye, Dongzhi Jiang, Zihao Wang, Leqi Zhu, Zhenghao Hu, Zilong Huang, Jun He, Zhiyuan Yan, Jinghua Yu, Hongsheng Li, Conghui He, and Weijia Li. Echo-4o: Harnessing the power of gpt-4o synthetic images for improved image generation, 2025. 5

[63] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Randomized autoregressive visual generation. *arxiv*, 2024. 2, 3

[64] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. In *International Conference on Learning Representations*, 2025. 2, 3

[65] 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. 4

[66] Xuanpu Zhang, Dan Song, Pengxin Zhan, Tianyu Chang, Jianhao Zeng, Qingguo Chen, Weihua Luo, and An-An Liu. Boow-vton: Boosting in-the-wild virtual try-on via mask-free pseudo data training. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 26399–26408. Computer Vision Foundation / IEEE, 2025. 9

[67] Boyang Zheng, Nanye Ma, Shengbang Tong, and Saining Xie. Diffusion transformers with representation autoencoders. *arXiv preprint arXiv:2510.11690*, 2025. 2, 3, 5, 8

[68] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1, 2, 4, 5, 8

[69] Lei Zhu, Fangyun Wei, Yanye Lu, and Dong Chen. Scaling the codebook size of vq-gan to 100,000 with a utilization rate of 99%. *Advances in Neural Information Processing Systems*, 37:12612–12635, 2024. 2

[70] Yongxin Zhu, Bocheng Li, Yifei Xin, Zhihua Xia, and Linli Xu. Addressing representation collapse in vector quantized models with one linear layer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22968–22977, 2025. 2