

# MARS: Mixture of Auto-Regressive Models for Fine-grained Text-to-image Synthesis

Wanggui He<sup>1,\*</sup>, Siming Fu<sup>1,\*</sup>, Mushui Liu<sup>2,\*</sup>, Xierui Wang<sup>2,+</sup>, Wenyi Xiao<sup>2,+</sup>, Fangxun Shu<sup>1,+</sup>, Yi Wang<sup>2</sup>, Lei Zhang<sup>2</sup>, Zhelun Yu<sup>3</sup>, Haoyuan Li<sup>2</sup>, Ziwei Huang<sup>2</sup>, LeiLei Gan<sup>2</sup>, Hao Jiang<sup>1,†</sup>,

<sup>1</sup>Alibaba Group <sup>2</sup>Zhejiang University <sup>3</sup>Fudan University

\*Equal contribution    †Core contributor    ‡Corresponding author

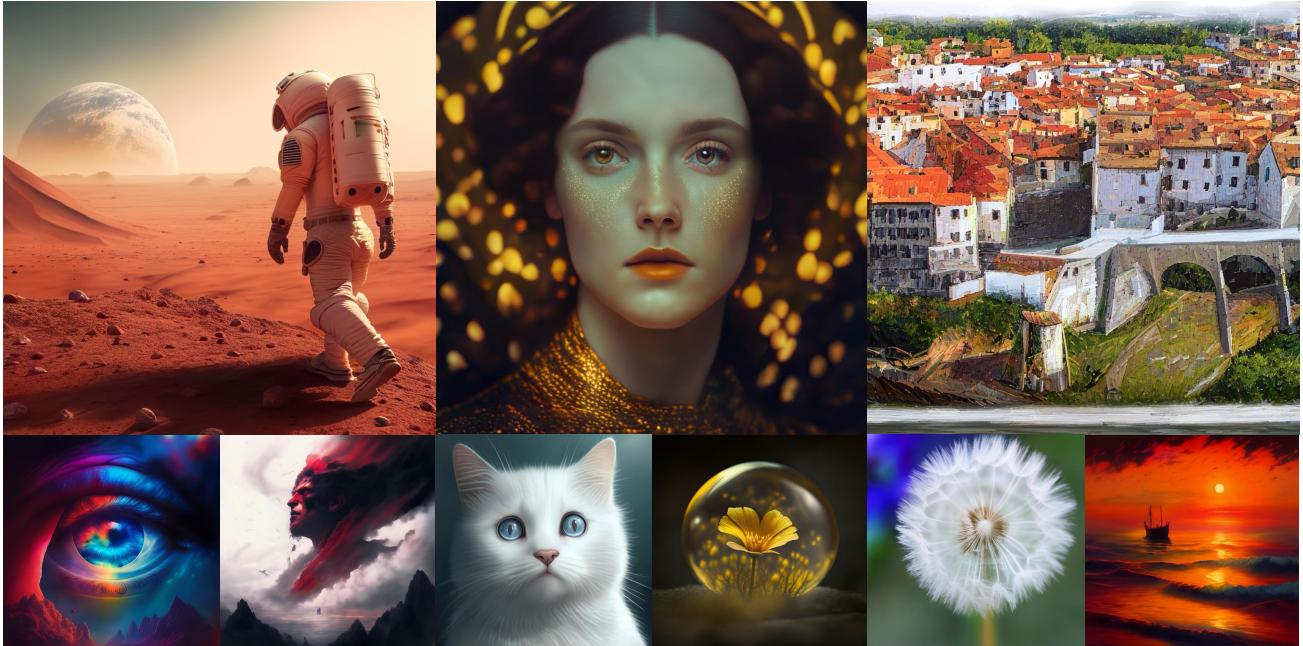


Figure 1. The generated samples from MARS display extraordinary quality, marked by an impressive degree of fidelity and precision in their adherence to the provided textual descriptions.

## Abstract

Auto-regressive models have made significant progress in the realm of language generation, yet do not perform on par with diffusion models in the domain of image synthesis. In this work, we introduce **MARS**, a novel framework for T2I generation that incorporates a specially designed Semantic Vision-Language Integration Expert (SemVIE). This innovative component integrates pre-trained LLMs by independently processing linguistic and visual information—freezing the textual component while fine-tuning the visual component. This methodology preserves the NLP capabilities of LLMs while imbuing them with exceptional visual understanding. Building upon the powerful base of the pre-trained Qwen-7B, MARS stands out with its bilingual generative capabilities corresponding to both English and

Chinese language prompts and the capacity for joint image and text generation. The flexibility of this framework lends itself to migration towards **any-to-any** task adaptability. Furthermore, MARS employs a multi-stage training strategy that first establishes robust image-text alignment through complementary bidirectional tasks and subsequently concentrates on refining the T2I generation process, significantly augmenting text-image synchrony and the granularity of image details. Notably, MARS requires only 9% of the GPU days needed by SD1.5, yet it achieves remarkable results across a variety of benchmarks, illustrating the training efficiency and the potential for swift deployment in various applications. Code will be available at <https://github.com/fusimeng3/MARS>.

## 1. Introduction

Pre-trained Large Language Models (LLMs) [3, 11, 57, 61, 66] have broadened their generative capabilities to encompass the visual domain. This advancement entails transforming pixel data into discrete tokens through a visual tokenizer, analogous to the processing of textual information, thereby integrating these tokens into the model’s transformer [59] architecture for generative tasks. Unlike other generative approaches, such as diffusion models [9, 17, 38, 45], LLMs [5, 14, 35, 64] uniquely utilize a discrete latent space of visual tokens, crucial for merging visual and linguistic modalities.

Auto-regressive models for text-to-image generation models, such as Parti [64], CogView2 [14], and Unified-io2 [34] have extended their generative scope to encompass the visual domain, facilitating the creation of images. These models integrate pre-trained LLMs within a unified architecture, enabling the simultaneous interpretation of both linguistic and visual inputs. Nonetheless, a notable challenge arises from the inherent distributional bias of LLMs, which are predominantly trained on textual data, potentially leading to a pronounced distributional shift when adapting to text-image pair datasets. This shift has the potential to provoke catastrophic forgetting, consequently impairing the LLMs’ primary competency in text generation tasks. *The aforementioned discourse prompts a pivotal inquiry: is it feasible to preserve the natural language processing proficiency of LLM while concurrently endowing it with state-of-the-art visual comprehension and generation capabilities?*

In response to this challenge, we present MARS, an innovative framework predicated on an auto-regressive model architecture akin to that of pre-trained LLMs for text-to-image synthesis. Specifically, we design the Semantic Vision-Language Integration Expert (SemVIE) module as the centerpiece of MARS to seamlessly facilitate the frozen pre-trained LLM with the trainable visual expert, thereby endowing them with exceptional visual understanding and preserving the NLP capability of pre-trained LLMs. Moreover, SemVIE can facilitate a comprehensive and incremental interplay between the textual and visual modalities across every layer of the model, fostering deep integration that yields images closely aligned with their textual descriptors. Through rigorous training on paired image-text datasets, MARS augments the generative capabilities of LLMs to include sophisticated text-to-image translations. As demonstrated in Fig. 1, MARS exhibits a pronounced ability to generate images with intricate visual details, such as animal fur, plant foliage, and facial features, underscoring its potent text-to-image generation proficiency.

In the domain of data optimization, we have developed a content-rich, efficient, and fine-grained approach for dataset construction. We leverage the capabilities of CogVLM [60] to generate sophisticated image descriptions that en-

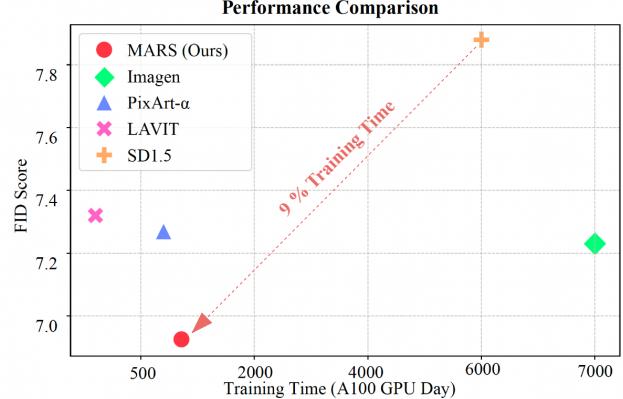


Figure 2. Comparison of training time and performance with models. The FID is evaluated on the zero-shot MS-COCO benchmark.

hance text-to-image alignment. For the optimization of the model training process, we have devised a multi-stage training strategy. This regimen begins with the creation of low-resolution images and advancing toward the production of high-resolution images with detailed textual alignment. Remarkably, with a mere 587 A100 GPU days, equating to only 9% of the training duration required by Stable Diffusion v1.5, MARS demonstrates its superiority over existing large-scale text-to-image (T2I) models, as evidenced in Fig. 2. Our contributions can be encapsulated as follows:

- We present MARS, an innovative framework adapted from auto-regressive pre-trained LLMs for T2I generation tasks. To ensure preservation of NLP capacities while also equipping the model with advanced visual generation and comprehension abilities, we design a module named SemVIE, which adds parallel visual experts to the attention blocks of pre-trained LLM. Therefore, MARS amplifies the flexibility of autoregressive methods for T2I generation and joint image-text synthesis, with the potential expansibility to **any-to-any** tasks.
- We propose a multi-stage refinement training strategy that significantly enhances MARS’ robust instruction-following capability and its ability to generate high-quality images with rich details.
- MARS shows great ability in prompt understanding and following, *e.g.* long and complex natural language inputs. Moreover, it possesses the **bilingual** capacity to follow prompts in both English and Chinese. The framework’s performance is verified across an array of evaluative measures, *i.e.* MS-COCO benchmark, T2I-CompBench, and Human Evaluation.

## 2. Related Works

### 2.1. Text-to-Image Generation Models

Text-to-image generation aims to create images based on given textual descriptions. Recent diffusion-based models

[24, 53–56] have demonstrated exceptional performance in image generation, offering improved stability and controllability. These models operate by introducing Gaussian noise to input images in a forward process and subsequently generate high-quality images with intricate details and diversity through an inverse process starting from random Gaussian noise. Models like GLIDE [36] and Imagen [48] utilize the CLIP [40] text encoder to enhance image-text alignment. Latent Diffusion Models (LDMs) [45] has been proposed to shift the diffusion process from pixel space to latent space, thereby enhancing efficiency and image quality. Furthermore, recent advancements such as SD-XL [38], DALL-E 3 [2], and Dreambooth [46] have significantly improved image quality and text-image alignment by employing various approaches, including innovative training strategies and scaling of training data. Furthermore, an architectural evolution is underway, with the diffusion model framework transitioning from a U-Net structure towards a transformer-based architecture DiT [37]. PixArt- $\alpha$  [9], SD-3.0 [17], and Lumina-T2X [21] achieve exceptional performance through the integration of DiT. The architecture evolution blurs the previously clear delineation between diffusion and language models in the visual generative arena. In this paper, we put forward a solution based on auto-regressive generation for better quality and interactive text-guided synthesis.

## 2.2. Auto-regressive Model for Visual Generation.

Auto-regressive Models [3, 11, 57, 61, 66] have been adeptly repurposed for the synthesis of visual media, including images [5, 14, 47] and videos [26, 62, 65]. The process begins with a visual tokenizer function implemented by VQ-VAE [58] or VQ-GAN [16],  $f$ , which effectively converts visual stimuli into a sequence of discrete tokens. Specifically, a video  $V \in \mathbb{R}^{T \times H \times W \times 3}$  (or an image when  $T = 1$ ) undergoes tokenization to yield a discrete representation  $X = f(V) \in \{1, \dots, K\}^{T_1 \times H_1 \times W_1}$ , where  $K$  denotes the codebook size intrinsic to the visual tokenizer [16]. Subsequently,  $X$  is linearized into a one-dimensional token sequence via raster scan order, which is then introduced to a language-model transformer to facilitate generative modeling. Current auto-regressive models, include notable architectures such as ImageGPT [10], DALL-E [42, 43], and Parti [63]. AR model anticipates the subsequent token based on a sequence of antecedent tokens, supplemented by additional conditional data  $c$ , and adheres to a categorical distribution for  $p_\theta(x_i | x_{<i}; c)$ .

## 3. Method

### 3.1. Preliminaries

**Auto-Regressive Models.** Auto-regressive models aim to predict future data points by regressing on their previous values. Current auto-regressive models are typically based

on Transformer-like architectures [59], leveraging the token prediction strategy.

**Next Token Prediction (NTP).** In the realm of sequential token analysis, one seeks to decipher the sequential arrangement, represented by the token sequence  $Z = \{z_1, z_2, \dots, z_{T_z}\}$ , in which each element  $z_t$  may correspond to either textual or pictorial information, encapsulated within a token, and  $T_z$  denotes the sequence's aggregate length. The endeavor of next token prediction (NTP) is directed towards the elucidation of the auto-regressive distribution  $P(z_{t+1} | z_{\leq t})$ , which characterizes the likelihood of each subsequent token, thus underpinning the generative process at each juncture of the sequence. The objective of NTP is elegantly quantified through Maximum Likelihood Estimation (MLE), harnessing the negative log-likelihood, equivalently appreciated as the cross-entropy loss, articulated mathematically as:

$$\mathcal{L}(\theta) = - \sum_{t=1}^{T_z-1} \log P_\theta(z_{t+1} | z_{\leq t}), \quad (1)$$

where  $\theta$  symbolizes the parameters that scaffold the model, with  $P_\theta(z_{t+1} | z_{\leq t})$  delineating the model's forecasted conditional probability distribution for the genesis of the ensuing token  $z_{t+1}$ .

**Next K-Token Prediction (NKTP)** Next Token Prediction offers the advantages of a straightforward task format and simplicity, as well as the ability to easily extend to text-image joint generation tasks. However, when generating high-resolution images, the requirement to output long sequences results in prolonged generation times and limited image quality. To address this issue, we propose utilizing Next K Token Prediction to enhance the resolution of images generated by Next Token Prediction. Specifically, NKTP extends the NTP framework by predicting the subset of the next  $K$  tokens instead of just the next single token based on predicted tokens. NKTP aims to capture longer dependencies and richer contextual information within the token sequence, enhancing the model's ability to generate coherent and contextually accurate sequences. In NKTP, the model learns to predict  $K$  tokens  $\{z^i, z^{i+1}, \dots, z^{i+K}\}$  given predicted tokens  $\{z^j | j \leq i\}$  at each auto-regressive step:

$$p(z_{i+1}, z_{i+2}, \dots, z_{i+K} | z_{\leq i}) \quad (2)$$

By considering multiple future tokens, NKTP can better model the dependencies between tokens, leading to more accurate and contextually appropriate predictions.

### 3.2. Overall Framework

We propose MARS, a confluence of large language models (LLMs) with vision generation capacities encapsulated within a unified framework. MARS embodies a balanced

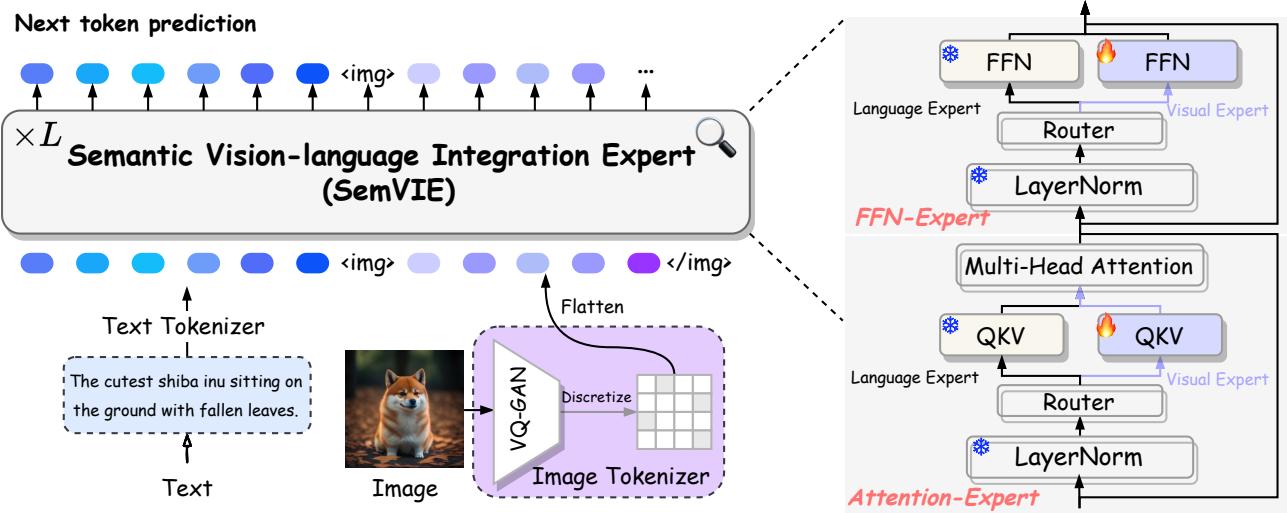


Figure 3. Overall training framework of the proposed MARS, which consists of the SemVIE modules facilitating T2I within a unified framework. An image-text pair is processed and tokenized by VQ-GAN [16] into ‘vision words’, which are then integrated with text tokens for joint processing in the SemVIE. The right part illustrates the multi-modal integration block, highlighting the synergistic processing of image and text data within the SemVIE, critical for the T2I task.

multi-modal architecture, comprising distinct yet harmonized visual and linguistic expert models, as delineated in Fig. 3. Consistency across modalities is sustained by parallel structural designs in both modules. The linguistic module leverages the capabilities of a pre-trained LLM, *e.g.* Qwen-7B [1], whereas the visual counterpart undergoes initialization concomitantly with the linguistic model. During the training phase, the linguistic component remains static, and optimization is confined to select weights within the visual domain, specifically calibrated for the image synthesis task. The architecture’s efficacy is further bolstered by an enriched visual vocabulary and the introduction of a *SemVIE*, which amalgamates the LLM’s sophisticated language interpretation abilities with visual perception. This cohesion not only harnesses the potent natural language processing capabilities inherent to the LLM but also supports the model’s education across a vast corpus of paired image-text exemplars, enhancing inter-modal congruity and fostering the generation of coherent visual content.

A detailed exposition of the *SemVIE* is outlined in Sec. 3.3. Subsequently, the manuscript explicates the nuanced process of multi-stage refinement in Sec. 3.4. We consummate the discussion with a presentation of the meticulously curated dataset of finely annotated image-text pairings in Sec. 3.5.

### 3.3. Semantic Vision-language Integration Expert

**Tokenization.** In this investigation, Qwen-7B [1], a pre-trained LLM serves as the foundational linguistic framework, leveraging its tokenizer to dissect the textual data

into a series of representative tokens denoted as  $r_t$ . Concurrently, within the visual modality, an encoder inspired by the VQ-GAN architecture [30] is employed to transform the image  $x \in \mathbb{R}^{3 \times H \times W}$  into a feature map  $f_v \in \mathbb{R}^{K \times D}$ , where,  $K = H \times W/P^2$ , with  $P$  predefined at a quantization parameter of 16, and  $D$  encapsulates the feature dimension. The feature map  $f_v$  is subsequently quantized using the visual codebook VQ-GAN that maps it onto a series of discrete code indices  $f_q$ . The process efficaciously refactors a  $256 \times 256$  pixel image into a sequence of 256 tokens, wherein each token embodies the information of a  $16 \times 16$ -pixel image segment. It is noteworthy that the visual codebook consists of 8192 unique codices. Such visual tokens are identified within the framework as  $r_v$ .

In the vocab of the MARS, these visual components are interwoven with traditional textual tokens, engendering a comprehensive multimodal vocabulary. The original vocab of the linguistic LLM encompasses 151,936 entries, which, upon symbiosis with the visual codebook and 6 special tokens (specifically designed to denote the start and end of image sequences, among other functionalities.), eventuates in a multimodal vocabulary size 160,136. Within the architecture of MARS, visual tokens synthesized by the VQ-GAN paradigm are conferred equitable status vis-à-vis their textual counterparts. Initial embeddings for the visual vocab are derived from the aggregative mean embedding of pre-trained textual tokens, establishing a foundational bedrock for ensuing cross-modality integrations.

**Semantic Vision-language Integration Expert.** The MARS architecture incorporates  $L$  layers of SemVIE,

which is a specialized multi-modal Mixture of Experts (mm-MoE) designed to adeptly handle both visual and semantic tokens. Central to the SemVIE are the Attention-MoE and Feed-Forward Network (FFN)-MoE modules. A dedicated routing module is strategically situated following each layer normalization step within the transformer modules. This routing mechanism is designed to allocate each input token to the corresponding expert model best equipped for its processing. A noteworthy aspect of the shared architectural framework is the universal application of the causal multi-head attention and layer normalization modules across both language and vision modalities, epitomizing a unified methodological approach to the concurrent processing of multi-modalities data. The process of Attention-MoE follows:

$$\begin{aligned} \hat{r}_t, \hat{r}_v &= \text{Router}(\text{LN}(\text{Concat}(r_t, r_v))) \\ \hat{r}_t^q, \hat{r}_t^k, \hat{r}_t^v &= W_Q^t(\hat{r}_t), W_K^t(\hat{r}_t), W_V^t(\hat{r}_t) \\ \hat{r}_v^q, \hat{r}_v^k, \hat{r}_v^v &= W_Q^v(\hat{r}_v), W_K^v(\hat{r}_v), W_V^v(\hat{r}_v) \\ \hat{r}_q, \hat{r}_k, \hat{r}_v &= \mathbf{C}(\hat{r}_t^q, \hat{r}_v^q), \mathbf{C}(\hat{r}_t^k, \hat{r}_v^k), \mathbf{C}(\hat{r}_t^v, \hat{r}_v^v) \\ \hat{r} &= \text{CausalAttention}(\hat{r}_q, \hat{r}_k, \hat{r}_v) + r \end{aligned} \quad (3)$$

where  $\mathbf{C}$  indicates concat operation,  $W_Q^t$ ,  $W_K^t$ , and  $W_V^t$  are frozen and loaded from pre-trained LLM.  $W_Q^v$ ,  $W_K^v$ , and  $W_V^v$  are trainable and initialized with the pre-trained semantic LLM. Then the MoE-FFN module further processes the multi-modal tokens:

$$\begin{aligned} \hat{r}_t, \hat{r}_v &= \text{Router}(\text{LN}(\mathbf{C}(r_t, r_v))) \\ \hat{r}_t &= \text{FFN}^t(r_t), \hat{r}_v = \text{FFN}^v(r_v) \\ \hat{r} &= \mathbf{C}(\hat{r}_t, \hat{r}_v) \end{aligned} \quad (4)$$

where  $\mathbf{C}$  indicates concat operation,  $\text{FFN}^t$  and  $\text{FFN}^v$  share the same architecture, and  $\text{FFN}^v$  is trainable. The SemVIE module, a cornerstone of the MARS, benefits from a synergistic integration of Attention-MoE and FFN-MoE modules, enabling the effective fusion of multimodal data streams. This integration capitalizes on the profound linguistic insights afforded by the pre-trained LLM, thus leveraging the advanced language comprehension capabilities to enrich visual understanding. To enable the model to simultaneously predict visual tokens and text tokens, in addition to using the original LLM model head (referred to as the text head), we added a vision head to the model. Notably, the text token and the visual token are processed through the text head and vision head to obtain the logits, denoted as  $l_t$  and  $l_v$ , respectively. The logits are then concatenated along the last dimension and passed through a softmax layer to obtain the probability distribution over the vocabulary for each token.

### 3.4. Multi-Stage Refinement

**Stage-I: Pre-training for Text-to-Image Alignment.** We first optimize MARS by two distinct tasks: text-to-image

generation and image captioning. This refinement process utilizes an auto-regressive approach for NTP, as explicated in Sec. 3.1. The procedure involves an extensive dataset of approximately 200 million text-image pairs, with each image conforming to a resolution of  $256 \times 256$  pixels.

**Stage-II: High-Quality Data Alignment.** To advance the fidelity of image synthesis, this stage persists in employing an NTP for the generation of images from textual descriptions. Diverging from Stage-I, the dataset enlisted for this stage comprises 50 million pairs of text and corresponding images, each pair meticulously curated through the application of an aesthetic valuation model [49]. The descriptive captions paired with these images originate from CogVLM [60], formulated in response to explicit directives. To mitigate potential discrepancies arising between the visual content and its textual descriptors, owing to image cropping, a standardized procedure is implemented wherein the minor axis of every image is resized to 256 pixels. This measure, taken whilst conserving the original aspect ratio, ensures the retention of comprehensive image content. However, this results in variable sequence lengths for the images. To address this, we include resolution information in the caption to specify the desired sequence lengths of the generated images.

**Stage-III: High-Resolution Refinement.** Inspired by the approaches of SD-XL [39] and DeepFloyd [13], we utilize a cascading super-resolution strategy to further enhance MARS. The low-resolution generated images and their corresponding captions serve as inputs to the super-resolution model. The super-res model is trained after the base model has been trained. In this stage, we employ the next K-token prediction (NTKP) method to predict higher-resolution images. The output images have a long side of **1024 pixels** while maintaining the original aspect ratio. To control the resolution of the generated images, we apply the same strategy as in Stage-II. Ten million triplet (low-resolution image, caption, high-resolution image) samples were used to train the cascaded super-resolution model.

### 3.5. Dataset Construction

The open-source English datasets incorporated into our study included LAION-400M [50], CC3M [52], CC12M [7], LAION-COCO [51], COYO [4], and Data-comp [19]. We initiate a filtration process to exclude images with resolutions below 256 pixels or aspect ratios greater than 2. Subsequently, we select images based on their CLIP scores [22] and aesthetic evaluations. This methodology yields a substantial corpus of 150 million image-text pairs. Additionally, we leveraged 50 million in-house data, predominantly comprising image-text pairs with Chinese captions, totaling approximately 200 million.

The coarse-grained image-text data exhibited substantial noise, evident in misalignments between images and

<b>Samples</b>			
<b>Original Caption</b>	A well-lit kitchen with granite countertops and bar stools surrounding an island.	A serene waterfall cascading in a lush forest with sunlight filtering through trees.	A spread of delicious Mediterranean dishes including falafel and various spices.
<b>Refined Caption</b>	The picture showcases a <b>luxurious</b> kitchen. The dominant color scheme is a combination of <b>cream</b> and <b>beige</b> , with <b>dark wooden</b> accents. The kitchen features a <b>large</b> central island with a granite countertop, set with <b>dishes</b> and a <b>bowl of bread</b> . There are <b>three</b> bar stools in front of the island. The cabinets are <b>tall</b> and have a <b>cream</b> finish, with some built-in appliances like an oven. The backsplash behind the stove is adorned with a <b>mosaic</b> pattern. The room is <b>well-lit</b> , with <b>natural</b> light coming through the island and <b>natural</b> light coming in from the windows. There are also decorative elements like a potted plant and a vase with bread on the countertop. <b>(78 words)</b>	The image showcases a <b>serene</b> forest setting with <b>tall, rugged</b> trees on the left. The trees have thick trunks and lush leaves, with moss growing on them. In the center, there's a cascading waterfall that flows smoothly, creating a misty effect. The waterfall is surrounded by dense foliage, predominantly in shades of <b>green</b> and <b>yellow</b> , indicating a transition between seasons. The lighting in the image is <b>soft</b> , possibly suggesting early morning or late afternoon, casting a <b>golden</b> hue on the scene. <b>(62 words)</b>	A <b>vibrant</b> and <b>appetizing</b> spread of Middle Eastern or Mediterranean food. There are <b>falafel balls</b> placed in a metal platter, with a few scattered on the wooden board below. Adjacent to the falafel, there's a plate of hummus garnished with chopped vegetables. A bowl of chickpeas in a spicy sauce is also visible. The table is adorned with various ingredients and accompaniments, such as branches of <b>fresh tomatoes</b> , a <b>lemon half</b> , <b>olives</b> in a bowl, and a glass of dark beverage. There are also some spices in small bowls, including what appears to be paprika and turmeric. The entire setup gives off a warm and inviting ambiance, making it ready to be enjoyed. <b>(75 words)</b>

Figure 4. Comparison of dataset captions before and after reconstruction by CogVLM [60]. The instruction prompt is *Describe the image and its style in a very detailed manner*. The **adjectives** are marked in blue, and **quantifiers** are marked in red to demonstrate the granularity of the reconstructed captions.

text, deficient descriptive content, irrelevant captions, and inferior image quality. To address these challenges in the succeeding T2I instruction following the training stage, we enhance the textual relevance and informational density through a caption rewriting strategy. Specifically, we deploy a pre-trained multimodal caption model CogVLM [60] to regenerate fine-grained captions for a curated selection of images. These newly generated captions intricately detail various aspects of the images, including object positioning, attributes, context, and stylistic elements, averaging approximately 110 words in length. Fig. 4 showcases an illustrative sample. This approach facilitated the generation of fine-grained captions for 50 million images.

## 4. Experiment

### 4.1. Experiment Details

**Implementation Details.** We employ AdamW [33] as the optimizer, with a beta parameter of 0.95 and weight decay set at 0.1. The peak learning rate is established at 1e-4, and a warm-up strategy is employed with a ratio of 0.01. For images with a resolution of  $256 \times 256$  pixels, the batch size per GPU is set at 64, while for  $512 \times 512$  pixel images, it is set at 24, leading to total batch sizes of 4096 and 1536, respectively. The training utilized DeepSpeed’s ZeRO-3 [41] optimization. The training epochs for Stage-I, Stage-II, and Stage-III of the model are configured to 1, 2, and 1 epochs, respectively.

**Evaluation Benchmarks.** We select three benchmarks for

comparison, including:

- **MSCOCO Dataset** [31]. Following previous works [15, 63], we generate 30k images use captions drawn from the MSCOCO 2014 evaluation dataset and assess both sample quality and image-text alignment of generated images. Specifically, we do not involve the selective curation of images from the generated output. The Fréchet Inception Distance (FID) [23] and CLIP Score [40] are used for evaluation.
- **T2I-CompBench** [27]. We employ various compositional prompts to assess textual attributes, including aspects such as color, shape, and texture, as well as attribute binding.
- **User Study.** We randomly select 100 prompts for evaluation. Subsequently, we enlist 30 participants for the user study.

### 4.2. Performance Comparisons and Analysis

**MSCOCO Benchmark.** We use the Frechet Inception Distance (FID) to evaluate the quality of synthesized images. As shown in 1, our proposed MARS, with only 7B trainable parameters, scores 6.92 on FID, which is a notable achievement. Compared to the auto-regressive counterpart Parti, we use fewer parameters (14B vs 20B) and smaller data sizes (0.2B vs 4.8B), achieving competitive performance (6.92 vs 7.22). Against the diffusion model SDv1.5, we achieve superior performance (6.92 vs 9.22) with less training budget (587 vs 6250 A100 GPU Days). These results highlight the efficiency of our mixture of auto-regressive

Table 1. Quantitative evaluation of FID and CLIPScore (where available) on MS-COCO 2014 for  $256 \times 256$  image resolution. **Diff** means diffusion model, **AR** means auto-regressive model. The results are all from the public literature. \* denotes that the results are picked from the different generated images with the best CLIP score.

Method	Venues	Architecture	#Params	FID-30K ↓	CLIPScore ↑
GLIDE [36]	ICML'22	Diff	5.0B	12.24	-
Imagen [25]	arXiv'22	Diff	3.4B	7.27	-
SDv1.0 [45]	CVPR'22	Diff	1B	-	30.50
SDv1.5 [45]	CVPR'22	Diff	-	9.22	-
MUSE [6]	ICML'23	Non-AR	3B	7.88	<u>32.00</u>
DALL-E 2 [42]	arXiv'22	Diff	3.5B	10.39	31.40
PixArt- $\alpha$ [9]	ICLR'24	Diff	7.32	-	
DALL-E [43]	ICML'21	AR	12.0B	28.00	-
CogView [15]	NeurIPS'21	AR	4.0B	27.10	-
Make-A-Scene [20]	ECCV'22	AR	4.0B	11.84	-
Parti [63]	arXiv'22	AR	20B	7.23	-
GILL [29]	NeurIPS'23	AR	6.7B	12.20	-
Emu [12]	arXiv'23	AR	13B	11.70	-
CM3Leon [64]*	arXiv'23	AR	7B	4.88	-
LAVIT [28]	ICLR'24	AR	7B	7.40	-
UIO-2 <sub>XXL</sub> [34]	CVPR'24	AR	-	13.39	-
MARS (Ours)	-	AR	14B	6.92	32.33
MARS* (Ours)	-	AR	14B	<b>3.51</b>	<b>33.10</b>

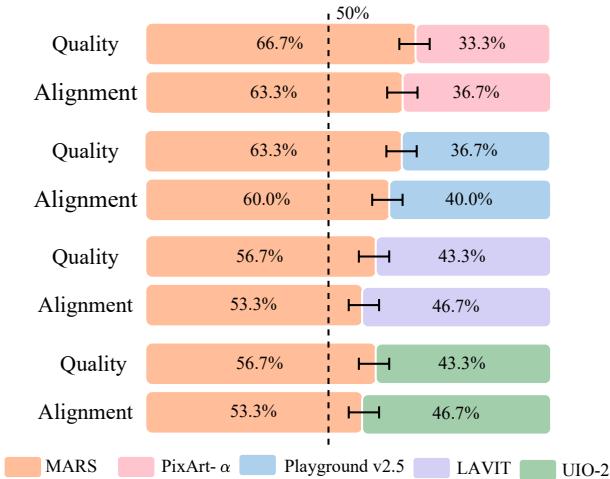


Figure 5. Human Evaluation Performance. Our MARS surpasses other state-of-the-art text-to-image models on both quality and alignment.

models.

Moreover, we utilize CLIP-Score to evaluate the alignment of textual conditions and corresponding generated images. MARS achieves 33.10 CLIPScore and 3.51 FID when the generated images are picked with the highest CLIP score, signaling its remarkable effectiveness in generating

visually compelling imagery that closely adheres to the semantic content of the text prompts.

**T2I CompBench Performance.** In the assessment of the T2I-CompBench, we curate a selection of contemporary text-to-image generative models for rigorous evaluation. This cohort includes Composable Diffusion [32], Structured Diffusion [18], Attn-Exct v2 [8], GORS [27], DALLE 2 [44], PixArt- $\alpha$  [9], SD1.5 [45], and SD-XL [38]. The empirical data presented in Tab. 2 delineates the superior performance of our proposed MARS within the T2I-CompBench benchmark, underscoring its proficiency in attribute binding, delineation of object relationships, and the synthesis of intricate compositions. Notably, MARS demonstrate a marked amelioration in the fidelity of color and texture representation, achieving enhancements of +11.63% in color fidelity and +7.49% in texture accuracy relative to DALL-E 2. It further exhibited substantial advancements in spatial and non-spatial metrics compared to DALL-E 2, with improvements quantified at +6.41% and +1.67%, respectively. Moreover, when juxtaposed with the recent PixArt- $\alpha$  model, which integrates a T5-XL text encoder, MARS outperforms it in various dimensions. Specifically, MARS achieved the highest scores in color (69.13%) and texture (71.23%) accuracy, outperforming PixArt- $\alpha$  which scored 68.86% and 70.44% respectively. These results demonstrate that the incorporation of LLM representations and visual tokens within an auto-

Table 2. Evaluation results (%) on T2I-CompBench [27]. The higher is better, and the best results are highlighted in bold.

Model	Venus	Attribute Binding			Object Relationship		Complex↑
		Color ↑	Shape ↑	Texture ↑	Spatial ↑	Non-Spatial ↑	
SD1.5 [45]	CVPR’22	37.65	35.76	41.56	12.46	30.79	30.80
SDXL [38]	arXiv’23	63.69	54.08	56.37	20.32	31.10	40.91
Composable Diffusion [32]	ECCV’22	40.63	32.99	36.45	8.00	29.80	28.98
Structured Diffusion [18]	ICLR’22	49.90	42.18	49.00	13.86	31.11	33.55
Attn-Exct v2 [8]	TOG’23	64.00	45.17	59.63	14.55	31.09	34.01
GORS [27]	ICCV’23	66.03	47.85	62.87	18.15	31.93	33.28
DALL-E 2 [44]	arXiv’22	57.50	54.64	63.74	12.83	30.43	36.96
PixArt- $\alpha$ [9]	ICLR’24	<u>68.86</u>	<b>55.82</b>	<u>70.44</u>	<b>20.82</b>	<u>31.79</u>	<b>41.17</b>
MARS (Ours)	-	<b>69.13</b>	<u>54.31</u>	<b>71.23</b>	<u>19.24</u>	<b>32.10</b>	<u>40.49</u>

Table 3. Ablation study of SemVIE on MS-COCO Benchmark. The term ‘w/o Visual Expert’ refers to a method wherein visual and text tokens are concatenated and used as inputs to fine-tune MARS without the implementation of the Visual Expert. Conversely, ‘w Visual Expert’ indicates the utilization of MARS’s specifically designed Visual Expert architecture.

Method	FID-30K ↓	CLIPScore ↑
w/o Visual Expert	10.13	30.14
w Visual Expert	8.24	31.03

Table 4. An ablation study assessing the impact of different stage training on the performance of the MARS model.

Method	FID-30K ↓	CLIPScore ↑
Stage-I	8.24	31.03
Stage-II	7.02	32.21
Stage-III	6.92	32.33

regressive framework can markedly improve the quality of generated images, as well as the alignment between the visual content and its corresponding textual narratives.

**User Study.** We conduct a user study evaluating various combinations of existing methods and MARS. Each combination is assessed based on two criteria: sample quality and image-text alignment. 60 Users are asked to evaluate the aesthetic appeal and semantic accuracy of images with identical text, determining which image is superior based on these criteria. Subsequently, we calculate the percentage scores for each model, as illustrated in Fig. 5. The results demonstrate that our MARS has significant advantages over both PixelArt- $\alpha$  and Playground-v2.5. Specifically, MARS achieves 66.7% and 63.3% higher voting preferences compared to PixelArt- $\alpha$  in terms of quality and alignment, respectively. Additionally, MARS shows a competitive per-

formance when compared to LAVIT and UIO-2.

### 4.3. Visual Analysis

Fig. 6 illustrates the sophisticated image synthesis capabilities of the MARS framework, producing visuals with remarkable detail and fidelity to textual descriptions. This proficiency is likely due to the advanced textual representations extracted from Large Language Models (LLMs), which, when integrated with a structured multi-tiered training strategy, significantly improve the model’s precision and alignment between text and image. The multi-stage training strategy of MARS incrementally refines the correlation between textual prompts and visual outputs, allowing for the generation of images that not only reflect the text’s intent but also display a depth of detail akin to photorealistic representations. Leveraging the deep semantic understanding from LLMs, MARS adeptly translates complex textual descriptions into coherent and contextually rich visual narratives, thus exemplifying a generative model that combines technical efficiency with artistic expression.

### 4.4. Multilingual Generation

Furthermore, at the heart of our language model lies the Qwen architecture, which is intrinsically designed to support multiple languages and incorporates a comprehensive dataset featuring both Chinese and English. During the training phase, a deliberate inclusion of a small yet significant proportion of Chinese in-house data. As depicted in Fig. 7, our model attains exemplary performance in Chinese text-to-image synthesis, notwithstanding the relative scarcity of Chinese corpus. This suggests that MARS has effectively mastered the ability to interpret concepts across linguistic boundaries, ensuring that both images and text coalesce within a singular representation space, as facilitated by our novel mixture mechanism.



A panoramic view of a city during the evening. The skyline is dominated by a mix of modern high-rise buildings and older architectural structures. The city is densely populated with buildings of varying heights and designs. In the foreground, there are residential buildings with balconies and trees. The city lights are on, casting a warm glow on the buildings.

Figure 6. Results of Visualization. The MARS framework is capable of generating realistic images across various resolutions and scenes.

#### 4.5. Ablation Study

We conduct ablation studies on the crucial parts discussed in Sec. 3.3 and Sec. 3.4, including model designs and multi-

stage training.

**Effect of SemVIE.** The results presented in Tab. 3 were obtained during Stage-I. The *w/o* Visual Expert configuration, which involves shared weights between the visual and lan-



这是一张黄昏时分的纽约城市景观照片，图中展现的是曼哈顿一侧的天际线，其中布鲁克林大桥的灯光在水面上留下了一道道光轨。地面上的车辆流光溢彩，仿佛将时间拉长一般。

图中是一坐高塔，这座高塔有七层，每层都有一些圆形的平面。塔的外观是棕色和白色。夕脚下，高塔附近树木的剪影显得很清晰。背景是蓝色的天空，没有云朵。下方的城市被树木和高塔遮挡，看起来像一座被绿色植物覆盖的山城。

这张图片展示了一杯冰咖啡，它装在一个塑料杯中，杯盖是球形的。杯子放在一张木制桌子上，桌面上有一个水杯，杯子的吸管是黑色的，放在杯子里。杯子的背景是一片绿色的草坪和一个花园。远处可以看到一些树木和一个模糊的建筑物。



图片展示了黄浦江边的城市风景，包括东方明珠电视塔、一些高楼大厦和一片蓝天。天空中飘着美丽的云彩，太阳可能即将落山。



图片展示的是一个美丽的游泳池。游泳池的水非常清澈，周围摆放着许多躺椅，白色遮阳伞随处可见，为游泳池边的休息区提供了阴凉。在游泳池的正面，有一个大喷泉，喷泉旁边是一片高尔夫球场，周围环绕着很多棕榈树。图片的最上方是蓝白相间的天空和满是绿植的山丘。

Figure 7. Illustration of the model’s multilingual capabilities. The model effectively responds to commands in Chinese, showcasing its inherent bilingual support and pointing toward future directions for multilingual research.



Figure 8. Qualitative comparison with different stages. The first row is from Stage II of MARS while the second row is from Stage III, indicating that Stage III is capable of correcting semantic inaccuracies while concurrently enhancing the quality of the images.

guage experts, leads to detrimental outcomes, evidenced by a 1.89 reduction in FID. This considerable decrease highlights the benefits of utilizing a specialized visual expert. The challenges associated with aligning visual and linguistic modalities underscore the need for specialized architectures that are adept at managing the intrinsic disparities between these types of data.

**Effect of Multi-Stage.** We further explore the effect of

training stages in Tab. 4. The results indicate that training in each stage positively impacts the model. On the MS-COCO benchmark, Stage II improved the FID by 1.22 compared to Stage I, and Stage III further enhanced it by 0.10 relative to Stage II. The visualizations of different stages are shown in Fig. 8. We observed that images generated during Stage I and Stage II lack detail, the images from Stage III exhibit the best quality.

为名为歌乐山辣子鸡的菜肴生成八幅步骤图像和相应的烹饪说明。描述：1. 木质表面展示各种食材：鸡腿、干辣椒、碎红辣椒、大蒜、生姜和葱。2. 切菜板上切有蔬菜片：姜片、根菜丝、大蒜片和葱。一只戴着塑料手套的手拿着一碗腌制的鸡肉。3. 旁边是切好的洋葱和绿色蔬菜。4. 炒锅里有沸腾的油和煎炸的食物。筷子将食物放在油上方。...

MARS:

**Text Output:** 菜名：“歌乐山辣子鸡”，“食材”：“鸡腿（300克）;盐: 4克;鸡精: 4克;糖: 3克;黄酒: 50克;山芋淀粉: 15克;花生油（烹饪鸡块用）: 500克;干辣椒: 100克;盐: 2克;白芝麻: 5克；，“步骤”：“准备食材。\\n鸡腿洗净剥骨后切成2厘米见方的块。（鸡腿肉口感细滑，多汁味美）\\n将切好的鸡块放入碗中，加入盐、鸡精、糖、黄酒、山芋淀粉、花生油，充分拌匀，静置15分钟备用。\\n干辣椒切段，蒜姜切片，葱切细丝。\\n锅中放油烧热，至五成热时放入鸡块，炸至表面金黄。（不要马上翻动，待定型后用筷子轻轻拨散）\\n将炸好的鸡块沥油捞出。\\n锅中留少量底油，将炸鸡余油再次烧热，倒入蒜姜、大葱和干辣椒炒香。\\n加一勺白芝麻。\\n倒入鸡块。\\n加一小勺盐。\\n加一勺鸡精。\\n倒入少量糖提鲜。\\n放入适量香油。\\n翻炒均匀出锅。”“配图”：<multi\_images>...</multi\_images>

**Image Output:**



为名为炸猪排的菜肴生成八个步骤图像和相应的烹饪说明。描述：1. 图像展示了摆放在台面上的各种食物和配料。左边有一个碗，里面装着细小的颗粒状物质，可能是面粉。旁边有一个鸡蛋。右边有两个小碗。一个碗里装着粉状物质，可能是糖，另一个碗里装着白色颗粒状物质，可能是盐或其他类型的面粉。中间有一个透明玻璃碗，里面装着一块生肉，可能是鸡肉或猪肉。2. 图像展示了生肉，可能是猪肉，放在一个...

MARS:

**Text Output:** "菜名": "炸猪排", "食材": "猪里脊: 300克; 面粉: 50克; 淀粉: 30克; 鸡蛋: 1个; 盐: 适量; 胡椒粉: 适量; 番茄酱: 50克; 千岛酱: 50克", "步骤": "1. 准备食材。猪肉里脊洗净，切成厚度约1-2厘米的片状。\\n2. 用肉锤或刀背捶打松里脊肉，这样炸起来比较嫩。\\n3. 鸡蛋打散，把捶打好的肉片放进去，并加入盐和胡椒粉，搅拌均匀。\\n4. 再加入玉米淀粉和面粉，确保每片肉都均匀裹上。\\n5. 腌制5分钟左右，使调料充分渗入。\\n6. 准备锅里加适量的油，烧到七成热时，改小火。\\n7. 拿出腌制好的肉片，轻轻放入油锅中油炸。\\n8. 炸的时候要不断翻面，炸至两面金黄，充分熟透后，捞出沥油。\\n将菜品装盘，配上番茄酱，享用美味的炸猪排。\\n配图: <multi\_images>...</multi\_images>"

**Image Output:**



Figure 9. Example of Multimodal Recipe Generation from MARS. MARS is capable of simultaneously generating text and images. The examples of the two recipes illustrated above demonstrate that, given the recipe title and the accompanying image captions, MARS can output the recipe steps and their corresponding images in an end-to-end manner. These generated images exhibit strong relevance to the text and maintain consistency and logical coherence among themselves.

#### 4.6. Further Analysis

**Image and Text Joint Generation Capability.** MARS extends beyond text-to-image generation, supporting the simultaneous generation of text and images, such as generating multiple text and image outputs from text and image inputs, with a focus on the relevance, consistency, and coherence between the two modalities. Due to the preservation of LLM’s integrity during MARS’s pre-training phase, the system is well-positioned for tasks involving concurrent text-image creation. For instance, in the domain of recipe generation, leveraging our text-image pre-trained model, we fine-tune it with a dataset of 10,000 recipes. This enables the model to produce comprehensive cooking tutorials that include step-by-step instructions accompanied by corresponding illustrations. As depicted in Fig. 9, upon receiving the recipe title and associated captions requiring images, the model concurrently generates detailed textual content, such as ingredient lists and procedural steps, as well as visual representations for each stage. Notably, MARS’s ability to seamlessly fuse text and imagery into coherent outputs are not confined to recipe generation and can be ex-

trapolated to other domains requiring joint text and image generation tasks.

#### 5. Conclusion

This study presents MARS, an innovative auto-regressive framework that not only retains the capabilities of pre-trained Large Language Models (LLMs) but also incorporates top-tier text-to-image (T2I) generation proficiency. MARS has been trained to exhibit exemplary performance in T2I tasks. We introduce the Semantic Vision-Language Integration Expert (SemVIE) module, which stands as the linchpin of MARS, streamlining the fusion of textual and visual token spaces and bringing a new insight into multimodal learning. MARS has demonstrated superior performance in multiple benchmark assessments, such as the MS-COCO benchmark, T2I-CompBench, and human evaluations. The pre-trained Qwen model equips MARS with the ability to generate bilingual images, blending Chinese and English seamlessly. Moreover, MARS adeptly handles joint image-text generation tasks, indicating its potential for any-to-any paradigm applications.

## References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 4
- [2] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023. 3
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020. 2, 3
- [4] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022. 5
- [5] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, José Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, Yuanzhen Li, Dilip Krishnan, and Google Research. Muse: Text-to-image generation via masked generative transformers. 2, 3
- [6] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, José Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. In *ICML*, 2023. 7
- [7] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. 5
- [8] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM TOG*, 42(4):1–10, 2023. 7, 8
- [9] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 2, 3, 7, 8
- [10] Mark Chen, Alec Radford, Rewon Child, Jeff Wu, Heewoo Jun, Prafulla Dhariwal, David Luan, and Ilya Sutskever. Generative pretraining from pixels. 2020. 3
- [11] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhang-hao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yong-hao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, 2023. 2, 3
- [12] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. 7
- [13] Deepfloyd. Deepfloyd. URL <https://www.deepfloyd.ai/>, 2023. 5
- [14] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. 2, 3
- [15] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. CogView: Mastering text-to-image generation via transformers. 2021. 6, 7
- [16] 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. 3, 4
- [17] 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. 2, 3
- [18] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. In *ICLR*, 2022. 7, 8
- [19] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacom: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36, 2024. 5
- [20] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 7
- [21] Peng Gao, Le Zhuo, Ziyi Lin, Chris Liu, Junsong Chen, Ruoyi Du, Enze Xie, Xu Luo, Longtian Qiu, Yuhang Zhang, et al. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers. *arXiv preprint arXiv:2405.05945*, 2024. 3
- [22] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 5
- [23] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. 2017. 6
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. 2020. 3
- [25] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Grishchenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 7
- [26] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv:2205.15868*, 2022. 3

- [27] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xi-hui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. In *ICCV*, 2023. 6, 7, 8
- [28] Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Chao Liao, Jian-chao Tan, Yadong Mu, et al. Unified language-vision pre-training in lilm with dynamic discrete visual tokenization. *arXiv preprint arXiv:2309.04669*, 2023. 7
- [29] Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. Generating images with multimodal language models. *NeurIPS*, 2023. 7
- [30] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image generation using residual quantization. In *CVPR*, 2022. 4
- [31] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 6
- [32] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *ECCV*, pages 423–439, 2022. 7, 8
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [34] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26439–26455, 2024. 2, 7
- [35] Xiaoxiao Ma, Mohan Zhou, Tao Liang, Yalong Bai, Tiejun Zhao, Huaian Chen, and Yi Jin. Star: Scale-wise text-to-image generation via auto-regressive representations. *arXiv preprint arXiv:2406.10797*, 2024. 2
- [36] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022. 3, 7
- [37] 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. 3
- [38] 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, 7, 8
- [39] 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. 5
- [40] 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 *ICML*, 2021. 3, 6
- [41] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models, 2020. 6
- [42] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. 3, 7
- [43] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3, 7
- [44] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv:2204.06125*, 2022. 7, 8
- [45] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 2, 3, 7, 8
- [46] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, pages 22500–22510, 2023. 3
- [47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamvar, Seyed Ghasemipour, Burcu Karagol, SSara Mahdavi, RaphaGontijo Lopes, Tim Salimans, Jonathan Ho, DavidJ Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. 3
- [48] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 3
- [49] Christoph Schuhmann. Clip+mlp aesthetic score predictor. <https://github.com/christophschuhmann/improved-aesthetic-predictor>. 5
- [50] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 5
- [51] C Schuhmann, A Köpf, R Vencu, T Coombes, and R Beaumont. Laion coco: 600m synthetic captions from laion2b-en. URL <https://laion.ai/blog/laion-coco>, 2022. 5
- [52] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018. 5
- [53] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265, 2015. 3
- [54] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. 2019.
- [55] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. 2020.

- [56] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 3
- [57] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2, 3
- [58] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *NeurIPS*, 2017. 3
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017. 2, 3
- [60] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023. 2, 5, 6
- [61] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *ICLR*, 2022. 2, 3
- [62] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 3
- [63] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 3, 6, 7
- [64] Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. *arXiv:2309.02591*, 2023. 2, 7
- [65] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 3
- [66] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. 2, 3