

OmniMamba: Efficient and Unified Multimodal Understanding and Generation via State Space Models

Jialv Zou^{1,◇} Bencheng Liao^{2,1,◇} Qian Zhang³ Wenyu Liu¹ Xinggang Wang^{1,✉}

¹ School of EIC, Huazhong University of Science & Technology

² Institute of Artificial Intelligence, Huazhong University of Science & Technology

³ Horizon Robotics

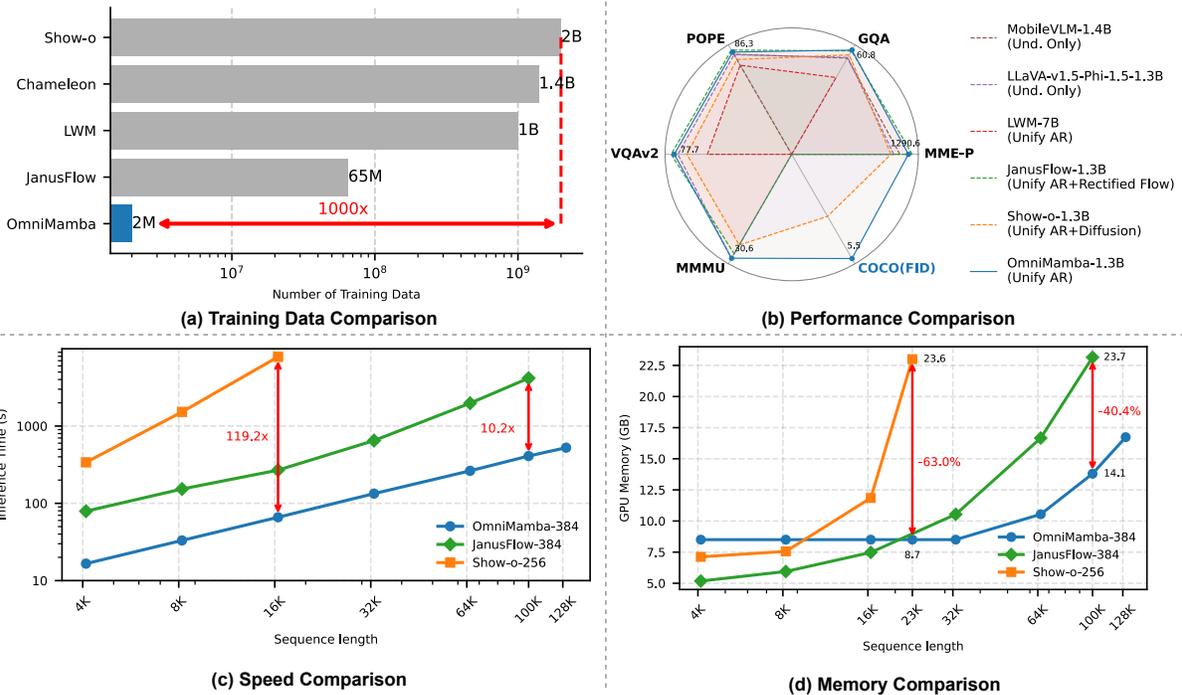


Figure 1. **Comprehensive comparison between OmniMamba and other unified understanding and generation models.** (a) Our OmniMamba is trained on only 2M image-text pairs, which is 1000 times less than Show-o. (b) With such limited data for training, our OmniMamba significantly outperforms Show-o across a wide range of benchmarks and achieves competitive performance with JanusFlow. Black metrics are for the multimodal understanding benchmark, while the blue metric is for the visual generation task. (c)-(d) We compare the speed and memory of OmniMamba with other unified models on the same single NVIDIA 4090 GPU. OmniMamba demonstrates up to a 119.2 \times speedup and 63% GPU memory reduction for long-sequence generation.

Abstract

Recent advancements in unified multimodal understanding and visual generation (or multimodal generation) models have been hindered by their quadratic computational complexity and dependence on large-scale training data. We present OmniMamba, the first linear-architecture-based multimodal generation model that generates both text and images through a unified next-token prediction paradigm.

[◇] Intern of Horizon Robotics.

[✉] Corresponding author: xgwang@hust.edu.cn

The model fully leverages Mamba-2’s high computational and memory efficiency, extending its capabilities from text generation to multimodal generation. To address the data inefficiency of existing unified models, we propose two key innovations: (1) decoupled vocabularies to guide modality-specific generation, and (2) task-specific LoRA for parameter-efficient adaptation. Furthermore, we introduce a decoupled two-stage training strategy to mitigate data imbalance between two tasks. Equipped with these techniques, OmniMamba achieves competitive performance with JanusFlow while surpassing Show-o across

benchmarks, despite being trained on merely 2M image-text pairs, which is 1,000 times fewer than Show-o. Notably, OmniMamba stands out with outstanding inference efficiency, achieving up to a 119.2× speedup and 63% GPU memory reduction for long-sequence generation compared to Transformer-based counterparts. Code and models are released at <https://github.com/hustvl/OmniMamba>

1. Introduction

In recent years, Large Language Models (LLMs) [2, 5, 15, 59, 60] have achieved remarkable advancements, igniting significant research interest in extending their fundamental capabilities to the visual domain. Consequently, researchers have developed a series of Multimodal Large Language Models (MLLMs) for tasks such as multimodal understanding [42, 43, 75, 77] and visual generation [31, 55].

Recent studies have emerged that seek to integrate multimodal understanding with visual generation, aiming to develop unified systems capable of handling both tasks simultaneously. Such designs hold the potential to foster mutual enhancement between generation and understanding, offering a promising pathway toward truly unifying all modalities. Numerous studies have sought to preserve the text generation paradigm of LLMs while exploring the impact [46, 64, 66, 67] of integrating diverse visual generation paradigms, such as diffusion models [24], flow-based generative models [16, 40], and vector-quantized autoregressive models [56].

Unfortunately, the significant domain gap between image and text presents a critical challenge for unified multimodal generative models: preserving generation capabilities without degrading understanding performance requires an extensive volume of image-text pairs for training, as illustrated in Fig. 1. This not only leads to poor training efficiency but also creates a substantial barrier to the broader development of such models, as only a small fraction of researchers possess the resources to undertake such computationally demanding studies. Moreover, most existing unified multimodal generative models rely on Transformer-based LLMs [61]. However, their quadratic computational complexity results in slow inference speeds, rendering them less practical for real-time applications.

The challenges faced by existing unified multimodal generative models naturally lead us to ponder: **can a model be developed that achieves both training efficiency and inference efficiency?**

To address this, we introduce OmniMamba, a novel unified multimodal generative model that requires only 2M image-text pairs for training. Built on the Mamba-2-1.3B [10] model as the foundational LLM with a unified next token prediction paradigm to generate all modalities,

OmniMamba leverages the linear computational complexity of state space models (SSMs) to achieve significantly faster inference speeds. Furthermore, to empower the Mamba-2 LLM—whose foundational capabilities are relatively weaker compared to the extensively studied Transformer models—to efficiently learn mixed-modality generation with limited training data, we propose novel model architectures and training strategies.

To enhance the model’s capability in handling diverse tasks, we incorporate task-specific LoRA [25]. Specifically, within each Mamba-2 layer’s input linear projection, we introduce distinct LoRA modules for multimodal understanding and visual generation. During task execution, the features are modulated by both the linear projection and the corresponding task-specific LoRA, while the irrelevant LoRA components are deactivated. Furthermore, we propose the decoupled vocabularies to guide the model in generating the appropriate modality, which requires more data for the model to learn. On the data front, we further propose a novel two-stage decoupled training strategy to address the data imbalance between the two tasks, significantly improving training efficiency.

Trained on only 2M image-text pairs, our proposed OmniMamba outperforms Show-o [67] on multiple multimodal understanding benchmarks and also matches the performance of JanusFlow [46], which was introduced by DeepSeek AI. Moreover, it achieves the best visual generation performance on the MS-COCO dataset [39]. Notably, OmniMamba demonstrates a 119.2× speedup at a sequence length of 16k and a 63% GPU memory reduction at a sequence length of 23k, compared to Show-o. Furthermore, at a sequence length of 100k, it achieves a 10.2× speedup and 40.4% memory savings compared to JanusFlow. Our main contributions can be summarized as follows:

- We introduce OmniMamba, the first Mamba-based unified multimodal understanding and visual generation model to the best of our knowledge. By novelly adopting decoupled vocabularies and task-specific LoRA, OmniMamba achieves effective training and inference.
- We propose a novel decoupled two-stage training strategy to address the issue of data imbalance between tasks. With this strategy and our model design, OmniMamba achieves competitive performance using only 2M image-text pairs for training-up to 1,000 times fewer than previous SOTA models.
- Comprehensive experimental results show that OmniMamba achieves competitive or even superior performance across a wide range of vision-language benchmarks and MS-COCO generation benchmark, with significantly improved inference efficiency, achieving up to a 119.2× speedup and 63% GPU memory reduction for long-sequence generation on NVIDIA 4090 GPU.

2. Related Work

Multimodal Understanding The remarkable advancements in LLMs have catalyzed the development of Large Vision-Language Models (LVLMs). Some representative works, such as the LLaVA series [43, 79], BLIP series [35, 36], and MiniGPT-4 [77], have demonstrated strong multimodal understanding capabilities. These models align the features obtained from pretrained vision encoders with the feature space of LLMs through feature projectors, enabling pretrained LLMs to transfer their understanding and reasoning abilities to multimodal scenarios.

Visual Generation In recent years, diffusion models [24] have made remarkable progress, leading to models [12, 49, 54] with strong visual generation capabilities. Building on these advancements, flow-based generative models [16, 40] have achieved superior results with fewer sampling steps. Additionally, some works [56, 71] have successfully integrated autoregressive models into this domain, achieving notable performance.

Unified Understanding and Generation The remarkable advancements of LLMs in the fields of multimodal understanding and visual generation have naturally sparked researchers’ interest in training a single LLM for both tasks. Early works [13, 18–20] integrated pretrained diffusion modules as tools into LLMs, essentially forming a combination of two expert systems rather than utilizing a single LLM to perform both tasks. This approach results in a more complex model architecture and often leads to suboptimal outcomes. Show-o [67] integrates next-token prediction with discrete diffusion, enabling adaptive handling of mixed-modality inputs and outputs. Meanwhile, JanusFlow [46] merges autoregressive models with rectified flow, a cutting-edge technique in visual generation. In contrast, Emu3 [64] asserts that next-token prediction holds the greatest potential for achieving multi-modal generation, relying exclusively on this paradigm to manage both text and image generation tasks. Although these methods achieve outstanding performance, they are all based on Transformers, whose quadratic computational complexity presents a significant drawback, particularly when handling long-sequence generation tasks and high-resolution image generation. To address this challenge, we propose OmniMamba, which employs a unified next token prediction paradigm to generate both text and image modalities, aiming to extend the linear computational complexity of the linear models to the field of multimodal generation.

Linear Model In recent years, a series of linear-complexity models have emerged as strong competitors to Transformers. Mamba [21], a selective state space model,

has garnered widespread attention for its competitive performance and faster inference speeds compared to Transformers. Building on this, Mamba-2 [10], an enhanced version of Mamba, achieves performance on par with Transformers while being 2-8 times faster. Similarly, GLA [68] leverages gated linear attention to achieve linear complexity, maintaining competitive performance with Transformers.

The success of linear-complexity models in the field of natural language processing (NLP) has inspired its application in the visual domain, where it has been extensively studied in traditional image tasks [45, 78], multimodal understanding [27, 38, 50, 74], and visual generation [26, 34]. In this paper, we aim to design a unified multimodal understanding and visual generation model based on Mamba-2, which maintains competitive performance with Transformer-based unified models while offering significantly faster inference speeds.

3. Method

3.1. Overall Architecture

Our ultimate goal is to design a unified multimodal understanding and visual generation model that achieves both training and inference efficiency using only 2M image-text pairs for training. We believe the key to realizing this goal can be summarized in one word: **decoupling**. To this end, we propose OmniMamba, the architecture of which is illustrated in Fig. 2.

Success necessitates standing on the shoulders of giants. We observe Emu3 [64], an autoregressive-based model which employs vast amounts of data and 8 billion model parameters. Despite these advantages, its final performance remains suboptimal, falling short of JanusFlow [46], a hybrid generative paradigm-based model with significantly less data and fewer parameters. We argue that this discrepancy stems not from the inherent superiority of the hybrid generative paradigm but from Emu3’s tight coupling design, it uses the same vocabulary and encoder for all tasks and modalities. While this design aligns with the original intention of a unified model, it may lead to inefficient data utilization. In the following, we will introduce our model by focusing on the concept of decoupling.

3.2. Decoupling Encoders for the Two Tasks

Previous works have explored using a single vision encoder for both tasks. For example, Show-o [67] employs MAGVIT-v2 [70] to encode images into discrete tokens for both understanding and generation tasks. TransFusion [76] utilizes a shared U-Net or a linear encoder to map images into a continuous latent space for both tasks. Emu3 trains its vision encoder based on SBER-MoVQGAN5, enabling the encoding of video clips or images into discrete tokens.

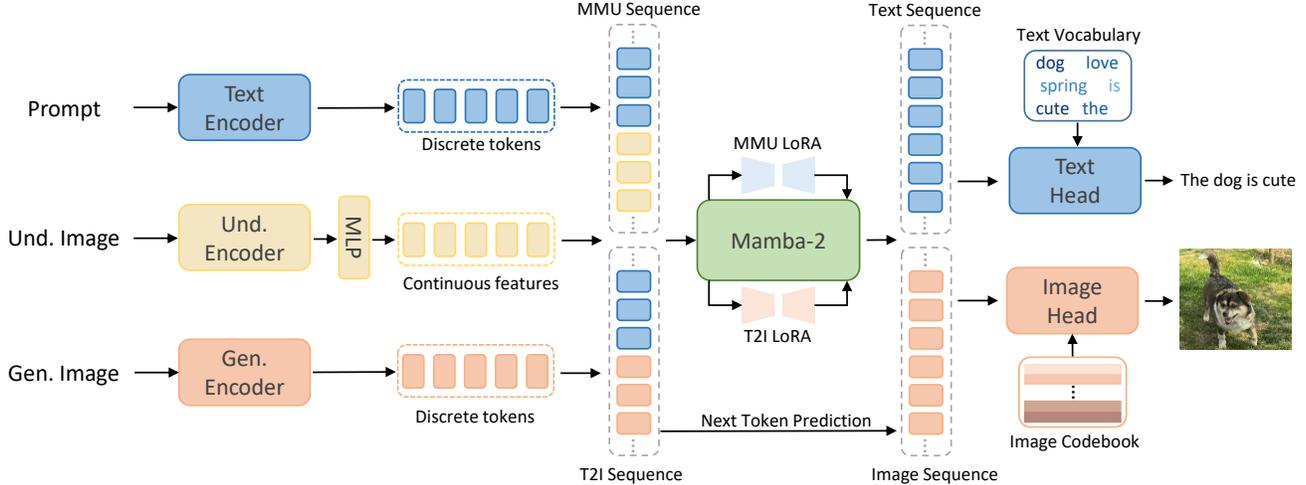


Figure 2. **Architecture of the proposed OmniMamba** “MMU” refers to multimodal understanding, while “T2I” refers to text-to-image generation. OmniMamba employs a next-token prediction paradigm for both multimodal understanding and visual generation tasks. To address the distinct requirements of each task—semantic information extraction for multimodal understanding and high-fidelity image compression for visual generation, we utilize separate encoders and heads. Furthermore, we purpose decoupled vocabularies to guide modality-specific generation and task-specific LoRA for parameter-efficient adaptation.

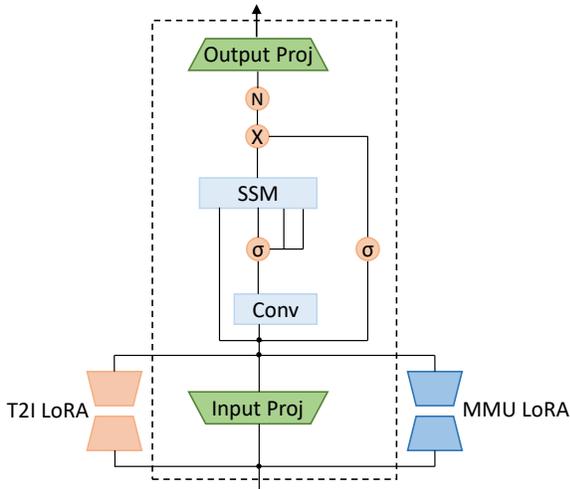


Figure 3. **The Mamba-2 block with task-specific LoRA.** It is worth noting that while the Mamba-2 Block in the Mamba-2 paper has two input projectors, the actual code implementation separates the feature dimensions from a single projector output. For simplicity, we depict only one input projector in our illustration. Our task-specific LoRA is applied to this entire input projector.

However, JanusFlow [46] has shown that such a unified encoder design is suboptimal. We believe this is primarily because multimodal understanding requires rich semantic representations for complex reasoning, whereas visual generation focuses on precisely encoding the spatial structure and texture of images. The inherent conflict between these two objectives suggests that a unified encoder design may not be the optimal choice. Therefore, OmniMamba adopts a decoupled vision encoder design.

Following prismatic VLMs [30], we fuse DINOv2 [48] and SigLIP [73] as an encoder to extract continuous features for multimodal understanding. The key idea is that integrating visual representations from DINOv2, which capture low-level spatial properties, with the semantic features provided by SigLIP leads to further performance improvements. For visual generation, we use an image tokenizer trained with LlamaGen [56] to encode images into discrete representations. This tokenizer was pretrained on ImageNet [11] and further fine-tuned on a combination of 50M LAION-COCO [33] and 10M internal high aesthetic quality data.

3.3. Decoupling Vocabularies for the Two Tasks

Unlike Emu3 and Show-o, which use a large unified vocabulary to represent both text and image modalities, to disentangle modality-specific semantics, we employ two separate vocabularies for each modality. This design explicitly separates the two modalities, providing additional modality-level prior knowledge. As a result, the model does not need to learn whether the output should be text or image, instead, it ensures the correct output modality by indexing the corresponding vocabulary. Our subsequent ablation experiments also confirm that OmniMamba’s dual-vocabulary design is one of the key factors for efficient training.

3.4. Task Specific LoRA

To enhance the model’s adaptability to specific tasks, we introduce task-specific adapters. We hypothesize that explicitly parameterizing the selection in SSMs based on task can enhance the data efficiency of multimodal training [14].

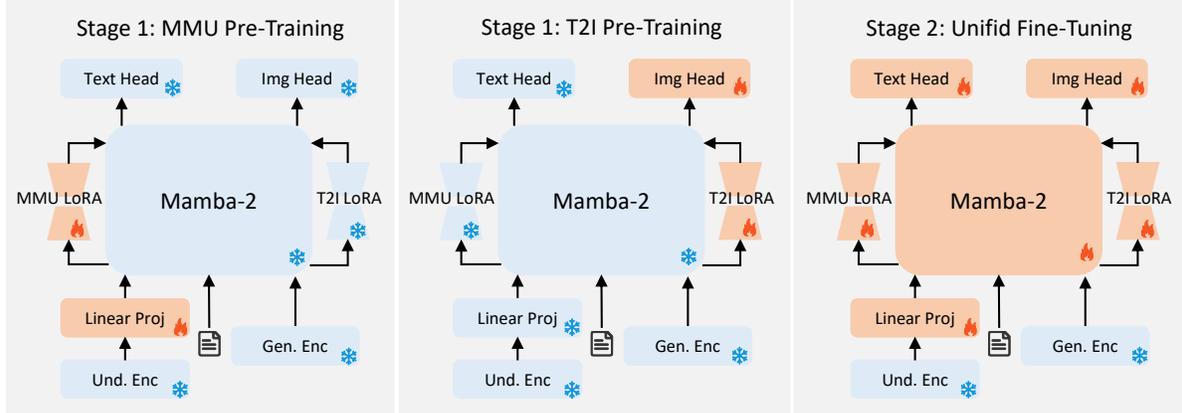


Figure 4. **Training strategy of OmniMamba.** The trainable components are indicated by a flame symbol, while the frozen ones are represented by snowflakes. The dashed arrows indicate that this route is temporarily dropped and does not participate in model training.

Specifically, to avoid introducing excessive parameters, we use LoRA [25] as the adapter. In OmniMamba, task-specific LoRA is applied only to the input projection of each Mamba-2 layer, as illustrated in Fig 3. When performing a specific task, the input linear projection and task-specific LoRA work together to effectively address the task. For instance, when the model performs a multimodal understanding (MMU) task, the MMU LoRA route is activated, while the text-to-image (T2I) LoRA route is dropped. Explicitly activating the corresponding adapter to assist in task execution helps improve data efficiency in training [14, 63].

3.5. Decoupled Training Strategy

We propose a decoupled two-stage training strategy to address data imbalance between understanding and generation tasks while improving training efficiency, as illustrated in Fig. 4. This approach consists of (1) a **Task-Specific Pre-Training** stage for module-specific initialization and modality alignment, and (2) a **Unified Fine-Tuning** stage for unified multi-task training.

Decoupling Rationale The first stage separates multimodal understanding (MMU) and text-to-image (T2I) generation tasks to prioritize modality alignment without data ratio constraints. Unlike joint pre-training methods (e.g., JanusFlow [46] with a fixed 50:50 MMU-T2I data ratio), our approach trains task-specific modules independently, enabling flexible dataset scaling (665K MMU vs. 83K T2I samples). Only randomly initialized components—linear projection and MMU LoRA for understanding, T2I LoRA and image head for generation—are trained, while the core Mamba-2 model remains frozen. This eliminates competition between tasks during early learning and allows asymmetric data utilization.

Stage 1: Task-Specific Pre-Training It contains: **MMU Pre-Training:** Trains the linear projection and MMU

LoRA to align visual-textual representations. The T2I LoRA path is disabled to isolate understanding-task learning. **T2I Pre-Training:** Optimizes the T2I LoRA and image decoder for visual synthesis. The MMU LoRA path is disabled to focus on generation capabilities.

Stage 2: Unified Fine-Tuning Inspired by multi-task frameworks [30, 79], we freeze the visual encoder and train all other modules while preserving task-specific LoRA independence. During each forward pass: (1) MMU and T2I computations use their respective LoRA branches; (2) Losses from both tasks are summed for a unified backward pass. This balances parameter sharing (via the frozen backbone) and task specialization (via isolated LoRA paths), enabling synergistic learning while mitigating interference between understanding and generation objectives.

3.6. Training Details

Data Formats Following Show-o [67], we use special tokens to unify the data formats for both multimodal understanding and visual generation tasks. The multimodal understanding data is structured as:

[MMU][SOI]{image tokens}[EOI][SOT]{text tokens}[EOT].

While the visual generation data is:

[T2I][SOT]{text tokens}[EOT][SOI]{image tokens}[EOI].

Specifically, [MMU] and [T2I] is a pre-defined task token used to guide the model in performing the corresponding task. [SOT] and [EOT] are used to represent the beginning and end of text tokens, respectively. Similarly, [SOI] and [EOI] represent the beginning and end of image tokens.

Training Objective Since OmniMamba uses the autoregressive paradigm to handle both multimodal understanding and visual generation tasks, we only need to use the

Type	Model	LLM Params	Res.	POPE↑	MME-P↑	VQAv2 _{test} ↑	GQA↑	MMMU↑
Und. Only	LLaVA-Phi [79]	Phi-2-2.7B	336	85.0	1335.1	71.4	-	-
	LLaVA [43]	Vicuna-7B	224	76.3	809.6	-	-	-
	Emu3-Chat [64]	8B from scratch	512	85.2	-	75.1	60.3	31.6
	LLaVA-v1.5 [42]	Vicuna-13B	448	86.3	1500.1	81.8	64.7	-
	InstructBLIP [8]	Vicuna-13B	224	78.9	1212.8	-	49.5	-
	MobileVLM [6]	MobileLLaMA-1.4B	336	84.5	1196.2	-	56.1	-
	MobileVLM-V2 [7]	MobileLLaMA-1.4B	336	84.3	1302.8	-	59.3	-
	LLaVA-v1.5-Phi-1.5 [67]	Phi-1.5-1.3B	336	84.1	1128.0	75.3	56.5	30.7
Unified	LWM [44]	LLaMA2-7B	256	75.2	-	55.8	44.8	-
	Chameleon [58]	7B from scratch	512	-	-	-	-	22.4
	LaVIT [29]	7B from scratch	256	-	-	66.0	46.8	-
	Emu3 [64]	8B from scratch	512	85.2	1243.8	75.1	60.3	31.6
	Janus [66]	DeepSeek-LLM-1.3B	384	87.0	1338.0	77.3	59.1	30.5
	JanusFlow [46]	DeepSeek-LLM-1.3B	384	88.0	1333.1	79.8	60.3	29.3
	Show-o [67]	Phi-1.5-1.3B	512	80.0	1097.2	69.4	58.0	26.7
	OmniMamba	Mamba-2-1.3B	384	86.3	1290.6	77.7	60.8	30.6

Table 1. **Comparison with other methods on multimodal understanding benchmarks.** “Und. only” refers to models that only perform multimodal understanding task, while “Unified” refers to models that unify both multimodal understanding and visual generation tasks. Models with a similar number of parameters to ours are highlighted in light blue for emphasis.

standard cross-entropy loss for next-token prediction during training.

4. Experiment

4.1. Data

To achieve the goal of data efficiency, we aim to train OmniMamba using as few high-quality image-text pairs as possible.

Multimodal Understanding Data In the first pretrain stage, the training data consists of 676K image-text pairs, all of which are sourced from publicly available datasets. These includes 118K images from COCO [39] and 558K images from LLaVA-1.5 pre-training data [42].

In the second fine-tune stage, we also exclusively use publicly available datasets follow Cobra [74].

1. The mixed dataset used in LLaVA-1.5 [42] consists of 665K visual multi-turn conversations.
2. LVIS-Instruct-4V [62], which comprising 220K images accompanied by visually aligned and context-aware instructions generated by GPT-4V.
3. LRV-Instruct [41], a large-scale visual instruction dataset of 400K samples, designed to mitigate hallucination issues across 16 vision-and-language tasks.

Visual Generation Data To facilitate better reproducibility and further exploration by the community, we using only 83K images from the MS-COCO 2014 dataset [39] for text-to-image generation training.

Overall, the training data for our unified multimodal generation model consists of fewer than 2 million image-text pairs. In contrast to previous works that rely on over 100M or even over 1B pairs, our approach is highly training efficient.

4.2. Implementation Details

Our core model is based on Mamba-2-1.3B, which consists of 48 layers of Mamba-2 blocks. In our primary experiments, the input image resolution for the multimodal understanding task is 384, while the image resolution for the visual generation task is 256.

For multimodal understanding, we combine DINOv2 [48] and SigLIP [73] as the image encoder, while for visual generation, we use the VQVAE trained by LlamaGen [56] as the image encoder. We incorporate task-specific LoRA into the input projector of each Mamba-2 block and set the LoRA rank to 8, which results in only a 0.65% increase in parameters. All training stages use the AdamW optimizer with β_1 set to 0.9 and β_2 set to 0.95. We adopt cosine annealing with warm-up as the learning rate schedule. Weight decay is set to 0, and gradient clipping is applied with a threshold of 1.0. Other detailed hyper-parameters are shown in the appendix. All of our training is conducted on NVIDIA A800 GPUs with BF16 precision.

4.3. Quantitative Results

Multimodal Understanding We evaluate OmniMamba’s multimodal understanding capabilities on a wide range of vision-language benchmarks, including POPE [37],

Type	Model	Params	Images	FID-30K↓
Gen. Only	DALL-E [53]	12B	250M	27.5
	GLIDE [47]	5B	250M	12.24
	DALL-E 2 [52]	6.5B	650M	10.39
	SDv1.5 [54]	0.9B	2000M	9.62
	PixArt [3]	0.6B	25M	7.32
	Imagen [55]	7B	960M	7.27
	Parti [69]	20B	4.8B	7.23
	Re-Imagen [4]	2.5B	50M	6.88
	U-ViT [1]	45M	83k(coco)	5.95
Unified	CoDI [57]	-	400M	22.26
	SEED-X [20]	17B	-	14.99
	LWM [44]	7B	-	12.68
	DreamLLM [13]	7B	-	8.76
	Show-o [67]	1.3B	35M	9.24
	OmniMamba	1.3B	83k(coco)	5.50

Table 2. **Compare visual generation capability with other methods on MS-COCO validation dataset.** “Gen. only” refers to models that only perform visual generation task, while “Unified” refers to models that unify both multimodal understanding and visual generation tasks.

Model	Gen_{avg} (Image/s)	Total (s)
Show-o [67]	0.81	19.66
JanusFlow [46]	1.02	15.64
OmniMamba	5.68	2.81

Table 3. **Image Generation Speed in Visual Generation Task.** OmniMamba achieves $7.0 \times$ faster image generation speed compared to Show-o and $5.6 \times$ faster compared to JanusFlow.

MME [17], GQA [28], MMMU [72]. The results are shown in Tab 1. Compared to models with a similar number of parameters, OmniMamba surpasses understanding-specific models such as LLaVA-v1.5-Phi-1.5 [67], MobileVLM [6], and MobileVLMv2 [7]. It also outperforms the unified understanding and generation model Show-o [67] and achieves competitive performance compare to the state-of-the-art unified model JanusFlow [46]. Notably, while Show-o utilizes 2B image-text pairs and JanusFlow leverages over 65M image-text pairs, OmniMamba achieves competitive performance by using only 2M image-text pairs for training.

Visual Generation We evaluate OmniMamba for text-to-image generation on the widely recognized MS-COCO [39] benchmark dataset. To quantify image quality, we report the FID score [23]. Consistent with previous literature, we randomly select 30K prompts from the MS-COCO validation set and generate corresponding images to compute the FID score. The results are shown in Tab 2, where our model

achieves the best visual generation performance on the MS-COCO validation dataset. Notably, models such as Show-o [67] and PixArt [69] are trained on external large-scale datasets and further fine-tuned on COCO-like datasets (e.g., OpenImages [32]) before evaluating zero-shot on the MS-COCO validation set. In contrast, to avoid introducing excessive additional data, both our model and U-ViT [1] are trained solely on the MS-COCO training set and evaluated on the MS-COCO validation set.

4.4. Qualitative Results

We present qualitative evaluations of our OmniMamba for both multimodal understanding and visual generation tasks. Fig 5 showcases our model’s capabilities in scene description and text-guided generation. Additional visualization results can be found in the appendix.

4.5. Inference Speed and GPU Memory usage

We compared the generation speed and GPU memory usage of OmniMamba with other Transformer-based models in both multimodal understanding and visual generation tasks. All the evaluations were done on the same single NVIDIA 4090 GPU with FP16 precision.

In multimodal understanding task, all models received the same example image. We used the same prompt, “Please describe the image in detail.” and removed the token generation limit to test their generation speed. The results are shown in the Fig 1. OmniMamba demonstrates $119.2 \times$ speedup at a sequence length of 16k, and saves 63.0% GPU memory at a sequence length of 23k compared to Show-o-256. Meanwhile, at a sequence length of 100k, OmniMamba achieves a $10.2 \times$ speedup and 40.4% GPU memory savings compared to JanusFlow-384, which is accelerated by FlashAttention-2 [9]. Notably, Show-o-256 indicates that the input image resolution is 256. Due to the design of its omni-attention mechanism being incompatible with FlashAttention-2, it was not used during testing. Similarly, JanusFlow-384 represents an input image resolution of 384, with FlashAttention-2 applied during testing for acceleration.

In visual generation task, we used the same prompt, “A Picture”. The models generate images with a batch size of 16 and a resolution of 256. As shown in Tab. 3, our model achieves image generation speeds that are $7.0 \times$ faster than Show-o and $5.6 \times$ faster than JanusFlow.

4.6. Ablation Studies

We conducted a series of ablation studies to verify the effectiveness of each design in OmniMamba. In this section, all ablation studies are conducted based on Mamba-2-370M, with the understanding visual encoder replaced by CLIP [51], an input resolution of 224, and a reduced number

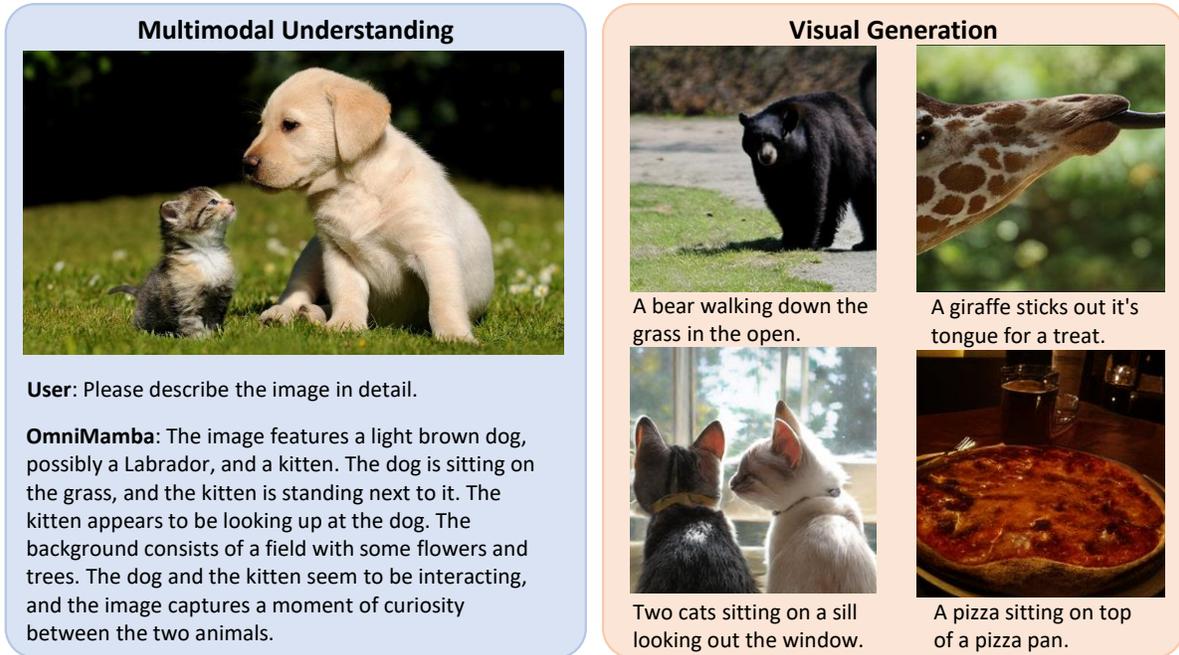


Figure 5. Qualitative results of OmniMamba on multimodal understanding and visual generation.

#Exp	Decoupling Vocabularies	Task Specific LoRA	POPE \uparrow	MME \uparrow	GQA \uparrow	FID-30K \downarrow
1	\times	\checkmark	80.8	1036	53.6	19.1
2	\checkmark	\times	81.2	1003	54.0	14.4
3	\checkmark	\checkmark	81.9	1100	55.3	10.3

Table 4. Ablation studies on decoupling Vocabularies and task specific LoRA in OmniMamba

of training steps (kept consistent across all ablation experiments), while keeping all other settings unchanged.

Impact of Decoupling Vocabulary for the Two Tasks

To guide the model in generating specific modalities from a structural design perspective, we employ modality-decoupled vocabularies in the default OmniMamba. We conducted ablation studies on this design. As shown in Tab 4, Exp1 and Exp3 utilize a modality-unified vocabulary and modality-decoupled vocabularies, respectively. The results demonstrate that the decoupled vocabularies enable more efficient model training and yield better performance. Notably, when using the modality-unified vocabulary for visual generation, the model occasionally produces text-related tokens, which requires additional post-processing to ensure correct visual generation. This further indicates that the model needs additional training to effectively learn modality-specific generation.

Impact of Task Specific Adapter We conducted ablation studies on the introduced task-specific adapter module, and the results are shown in Tab 4. Exp2 and Exp3 represent the

model without and with the task-specific adapter, respectively. The experiments demonstrate that the task-specific adapter helps the model efficiently learn both multimodal understanding and visual generation with a minimal amount of training image-text pair data (2M).

5. Conclusion

We presented OmniMamba, the first Mamba-2-based unified multimodal understanding and visual generation framework that achieves competitive performance with remarkable inference and training efficiency. By introducing three key innovations: decoupled vocabularies to disentangle modality-specific semantics, task-specific LoRA modules for parameter-efficient adaptation, and a two-stage decoupled training strategy to resolve data imbalance, OmniMamba achieves comparable performance with JanusFlow and even surpasses Show-o using only 2M image-text pairs for training. Moreover, OmniMamba exhibits outstanding inference efficiency, It achieves a $119.2\times$ speedup with a sequence length of 16k and a 63% reduction in GPU memory at a sequence length of 23k, compared to Show-o. With

a sequence length of 100k, it delivers a $10.2\times$ speedup and saves 40.4% of GPU memory compared to JanusFlow. These results validate that our proposed OmniMamba is both training and inference efficient, with the potential to enable more ordinary researchers to participate in the wave of unified model innovation. However, due to the limited scale of training data, our model’s performance remains slightly below SOTA methods. Exploring the trade-off between training data volume and model performance will be a key focus of our future work.

References

- [1] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22669–22679, 2023. 7
- [2] 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. *Advances in neural information processing systems*, 33:1877–1901, 2020. 2
- [3] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 7
- [4] Wenhui Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *arXiv preprint arXiv:2209.14491*, 2022. 7
- [5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023. 2
- [6] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023. 6, 7
- [7] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*, 2024. 6, 7
- [8] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023. 6
- [9] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023. 7
- [10] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024. 2, 3
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3
- [13] Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. *arXiv preprint arXiv:2309.11499*, 2023. 3, 7
- [14] Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, et al. Loramoe: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment. *arXiv preprint arXiv:2312.09979*, 4(7), 2023. 4, 5
- [15] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahri, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 2
- [16] 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
- [17] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. 7
- [18] Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. Planting a seed of vision in large language model. *arXiv preprint arXiv:2307.08041*, 2023. 3
- [19] Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer. *arXiv preprint arXiv:2310.01218*, 2023.
- [20] 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. 3, 7
- [21] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 3
- [22] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 12
- [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. *Advances in neural information processing systems*, 30, 2017. 7
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information*

- processing systems*, 33:6840–6851, 2020. 2, 3
- [25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 5
- [26] Vincent Tao Hu, Stefan Andreas Baumann, Ming Gui, Olga Grebenkova, Pingchuan Ma, Johannes Fischer, and Björn Ommer. Zigma: A dit-style zigzag mamba diffusion model. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024. 3
- [27] Wenjun Huang, Jiakai Pan, Jiahao Tang, Yanyu Ding, Yifei Xing, Yuhe Wang, Zhengzhuo Wang, and Jianguo Hu. Ml-mamba: Efficient multi-modal large language model utilizing mamba-2. *arXiv preprint arXiv:2407.19832*, 2024. 3
- [28] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 7
- [29] Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Quzhe Huang, Bin Chen, Chenyi Lei, An Liu, et al. Unified language-vision pretraining in llm with dynamic discrete visual tokenization. arxiv 2024. *arXiv preprint arXiv:2309.04669*. 6
- [30] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. *arXiv preprint arXiv:2402.07865*, 2024. 4, 5
- [31] Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal language models. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [32] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International journal of computer vision*, 128(7):1956–1981, 2020. 7
- [33] LAION. Laion-coco 600m. <https://laion.ai/blog/laion-coco>, 2022. 4
- [34] Haopeng Li, Jinyue Yang, Kexin Wang, Xuerui Qiu, Yuhong Chou, Xin Li, and Guoqi Li. Scalable autoregressive image generation with mamba. *arXiv preprint arXiv:2408.12245*, 2024. 3
- [35] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 3
- [36] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022. 3
- [37] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. 6
- [38] Bencheng Liao, Hongyuan Tao, Qian Zhang, Tianheng Cheng, Yingyue Li, Haoran Yin, Wenyu Liu, and Xing-gang Wang. Multimodal mamba: Decoder-only multimodal state space model via quadratic to linear distillation. *arXiv preprint arXiv:2502.13145*, 2025. 3
- [39] 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 *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2, 6, 7
- [40] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 2, 3
- [41] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning. In *The Twelfth International Conference on Learning Representations*, 2023. 6
- [42] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 2, 6
- [43] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2, 3, 6
- [44] Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv e-prints*, pages arXiv–2402, 2024. 6, 7
- [45] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin Jiao, and Yunfan Liu. Vmamba: Visual state space model, 2024. 3
- [46] Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan, Zhenda Xie, Haowei Zhang, Liang Zhao, et al. Janusflow: Harmonizing autoregression and rectified flow for unified multimodal understanding and generation. *arXiv preprint arXiv:2411.07975*, 2024. 2, 3, 4, 5, 6, 7
- [47] 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. *arXiv preprint arXiv:2112.10741*, 2021. 7
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 4, 6
- [49] 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. 3
- [50] Yanyuan Qiao, Zheng Yu, Longteng Guo, Sihan Chen, Zijia Zhao, Mingzhen Sun, Qi Wu, and Jing Liu. V1-mamba: Exploring state space models for multimodal learning. *arXiv preprint arXiv:2403.13600*, 2024. 3
- [51] 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. 7
- [52] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 7
- [53] 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. 7
- [54] 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*, pages 10684–10695, 2022. 3, 7
- [55] 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. *Advances in neural information processing systems*, 35:36479–36494, 2022. 2, 7
- [56] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 2, 3, 4, 6
- [57] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. *Advances in Neural Information Processing Systems*, 36, 2024. 7
- [58] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024. 6
- [59] 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
- [60] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2
- [61] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 2
- [62] Junke Wang, Lingchen Meng, Zejia Weng, Bo He, Zuxuan Wu, and Yu-Gang Jiang. To see is to believe: Prompting gpt-4v for better visual instruction tuning. *arXiv preprint arXiv:2311.07574*, 2023. 6
- [63] Peng Wang, Shijie Wang, Junyang Lin, Shuai Bai, Xiaohuan Zhou, Jingren Zhou, Xinggang Wang, and Chang Zhou. One-peace: Exploring one general representation model toward unlimited modalities. *arXiv preprint arXiv:2305.11172*, 2023. 5
- [64] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 2, 3, 6
- [65] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. 12
- [66] 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. 2, 6
- [67] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 2, 3, 5, 6, 7
- [68] Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023. 3
- [69] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gungjan 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*, 2(3):5, 2022. 7
- [70] Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion—tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023. 3
- [71] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Randomized autoregressive visual generation. *arXiv preprint arXiv:2411.00776*, 2024. 3
- [72] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 7
- [73] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023. 4, 6
- [74] Han Zhao, Min Zhang, Wei Zhao, Pengxiang Ding, Siteng Huang, and Donglin Wang. Cobra: Extending mamba to multi-modal large language model for efficient inference. *arXiv preprint arXiv:2403.14520*, 2024. 3, 6
- [75] Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. Tinyllava: A framework of small-scale large multimodal models. *arXiv preprint arXiv:2402.14289*, 2024. 2
- [76] Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024. 3

- [77] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 2, 3
- [78] Lianghai Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv preprint arXiv:2401.09417*, 2024. 3
- [79] Yichen Zhu, Minjie Zhu, Ning Liu, Zhiyuan Xu, and Yaxin Peng. Llava-phi: Efficient multi-modal assistant with small language model. In *Proceedings of the 1st International Workshop on Efficient Multimedia Computing under Limited*, pages 18–22, 2024. 3, 5, 6

Appendix

A. Training Details

The detailed training hyper-parameters are listed in Tab 5. The first stage separates multimodal understanding (MMU) and text-to-image (T2I) generation tasks to prioritize modality alignment without data ratio constraints. To this end, we employ a larger learning rate during the first pre-training stage and a smaller learning rate in the second fine-tuning stage. The batch size ratio represents the proportion between multimodal understanding data and visual generation data. All training is conducted on NVIDIA A800 GPUs using BF16 precision.

	Stage 1: MMU	Stage 1: T2I	Stage 2: Unify
Learning Rate	1e-3	8e-4	1e-4
Warm-up Steps	100	1000	0
Training Steps	5k	100k	150k
Batch Size	256:0	0:720	3:48

Table 5. **Hyper-parameters of OmniMamba.** The batch size ratio refers to the proportion between multimodal understanding data and visual generation data.

B. Limitations

Although our OmniMamba achieves promising results with a very small amount of data, the limited data volume still renders our model suboptimal. Furthermore, unlike previous works that leverage large-scale, high-quality datasets such as LAION-aesthetics, we rely solely on the MS-COCO dataset for visual generation. As a result, the quality of generated images, particularly for human faces, remains less refined. Exploring the trade-off between dataset scale and model performance will be a key focus of our future work.

Additionally, while Mamba-2 demonstrates exceptional inference efficiency, its foundational capabilities remain weaker compared to the extensively studied Transformer. Furthermore, Mamba-2 has only been trained on sequences

of up to 2048 tokens, limiting its ability to handle ultra-long sequences and hindering its extension to advanced techniques such as Chain-of-Thought (CoT) [65] or reinforcement learning [22]. Enhancing Mamba-2’s foundational capabilities and its capacity to model ultra-long sequences will be critical areas for future investigation.

C. Additional Qualitative Results

C.1. Multimodal Understanding

We validate the multimodal understanding capabilities of our OmniMamba and other approaches across three aspects: scene description, spatial reasoning, and counting with tricky questions, which are shown in Table 6–8.

Spatial Reasoning:



User	Is the bicycle parked to the right of the dog in this image?
OmniMamba	The bicycle is actually parked to the left of the dog, not to the right.
Show-o	Yes , the bicycle is parked to the left of the dog.
JanusFlow	No.

Table 6. OmniMamba can accurately determine spatial positions, JanusFlow made the correct judgment, but its response was overly concise, while Show-o lacks spatial reasoning capabilities, as indicated by the **red** text.

Counting with Tricky Questions:



User	The animals in the picture look like cows. How many are there?
OmniMamba	There are three goats in the image, not cows.
Show-o	There are three cows in the image.
JanusFlow	There are three animals in the picture.

Table 7. OmniMamba correctly outputs the quantity while avoiding the user’s trap question, whereas Show-o, despite counting correctly, was misled by the question, as indicated by the **red** text. On the other hand, JanusFlow’s response was not accurate enough.

Scene Description:



User

Please describe it in detail.

OmniMamba

The image depicts a serene scene of a cat and a dog sharing a peaceful moment together. The cat, exhibiting a gray and white coat, is curled up comfortably, its head resting on the dog's body. The dog, with its brown fur, is curled up in a similar position, with its head resting on the cat's back. Their eyes are gently closed, signifying a deep sleep. The background is a blurred white, allowing the focus to remain on the sleeping feline and its canine friend. The image exudes warmth and affection, showcasing the close bond between the two animals.

Show-o

The image features a white and **brown** cat lying on a bed, curled up and sleeping. The cat is resting its head on a **stuffed animal, which is a teddy bear. The teddy bear is positioned on the right side of the cat, with its head resting on the cat's back. The scene captures a peaceful and cozy moment between the cat and the stuffed animal.**

Table 8. OmniMamba can accurately describe the information in the scene, whereas Show-o made a mistake about the color of the cat and misidentified the dog as a teddy bear, as indicated by the **red** text.

C.2. Visual Generation

We provide additional visualization results of visual generation to further validate our generation capabilities, as shown in Fig 6.



A close up view of a guy brushing his teeth.



A colorful bird perched on a branch in the wild.



A close shot of a computer keyboard.



A laptop computer sitting on top of a desk.



A man cutting off his tongue with scissors.



A sink under a large mirror in a bathroom.



A female tennis player on a tennis court.



A good looking pizza is still in the box.



A surfer in a wet suit is surfing on a white board.



A jumbo jet plane flying through the air in a cloud filled sky.



A bathroom with a toilet, cabinet and rug.



A bathroom with two sinks mounted on a wall.



A bed in a bedroom between two lamps.



A bowl that has different types of food in it.



A brown table with white plate holding a pizza.

Figure 6. Qualitative results of OmniMamba visual generation. Prompts are randomly drawn from the MS-COCO validation set.