

# BEYOND TEXT-TO-IMAGE: LIBERATING GENERATION WITH A UNIFIED DISCRETE DIFFUSION MODEL

Anonymous authors

Paper under double-blind review

## ABSTRACT

Unified generation models aim to handle diverse tasks across modalities—such as text-to-image generation and image-to-text generation—within a single architecture and decoding paradigm. Autoregressive unified models suffer from slow inference due to sequential decoding, and non-autoregressive unified models suffer from weak generalization due to limited pretrained backbones. We introduce Muddit, a **unified discrete diffusion transformer** that enables fast and parallel generation across both text and image modalities. Unlike prior unified diffusion models trained from scratch, Muddit integrates strong **visual priors** from a pre-trained text-to-image backbone with a lightweight text decoder, enabling flexible and high-quality multimodal generation under a unified architecture. Empirical results show that Muddit achieves competitive or superior performance compared to significantly larger autoregressive models in both quality and efficiency. The work highlights the potential of purely discrete diffusion, when equipped with strong visual priors, as a scalable and effective backbone for unified generation.

## 1 INTRODUCTION

Unified generative models have recently emerged as a promising paradigm for multimodal data, encompassing both text and images. Most existing approaches adopt the autoregressive (AR) framework (Touvron et al., 2023), where modalities are represented as discrete token sequences and generated sequentially in raster order. While this paradigm is well-suited for language, it introduces severe inefficiencies in image generation: producing an image requires step-by-step prediction of thousands of tokens, leading to substantial computational cost. Moreover, the imposed rasterized order is poorly aligned with the inherently two-dimensional structure of images. These limitations hinder speed/quality trade-offs and restrict flexible conditional generation, such as inpainting, thereby constraining the practical applicability of unified models in interactive or real-time scenarios. To mitigate these issues, recent works (Chen et al., 2025a; Pan et al., 2025; Chen et al., 2025b) have proposed hybrid approaches that couple AR-based language models with diffusion-based image generators (Ho et al., 2020), as shown in Fig. 1 (a). However, such “glue” architectures fall short of true unification, as they introduce additional complexity into the inference pipeline while retaining considerable computational overhead. So there is a lack of a principled multimodal generative paradigm over current unified models.

As shown in Fig 1 (b), recent work like Dual-Diffusion (Li et al., 2024c) explores unifying multimodal under the diffusion model, but it ultimately relies on continuous diffusion for image (Esser et al., 2024) and discrete diffusion for text (Swerdlow et al., 2025b). This fundamental mismatch in generative principles undermines its claim of a true unification paradigm. UniDisc (Swerdlow et al., 2025a) takes a more promising step by applying discrete diffusion over multimodal token spaces<sup>1</sup>. This allows parallel refinement of text and image tokens, improving inference efficiency and enabling more flexible conditioning. However, the overall quality of UniDisc’s generation remains far from satisfactory. For example, it fails to match the fidelity of early diffusion models such

<sup>1</sup>MaskGIT, MaskAR, RandomAR, and Discrete Diffusion share significant conceptual and practical overlaps, often differing only in decoding order or architectural nuances. We elaborate on their connections in the next section. While Meissonic (Bai et al., 2025) follows the naming convention of MaskGIT (Chang et al., 2022), we standardize terminology in this paper by referring to all such models under the umbrella of Discrete Diffusion.

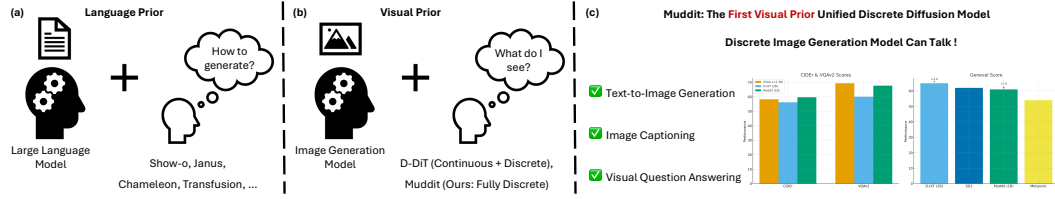


Figure 1: We propose Muddit, the first unified discrete diffusion model with a visual prior. Compared to language prior models like Show-o (Xie et al., 2024), Muddit demonstrates strong performance at image captioning and visual question answering. It also delivers clearer gains (7.0 vs 3.0) in image generation over the visual prior model D-DiT (Li et al., 2024c).

as Stable Diffusion 1.5 (Rombach et al., 2022), and lacks support for vision-language reasoning tasks such as visual question answering (VQA). We attribute these shortcomings to the pretrained model’s lack of prior knowledge. Without modular components carrying rich priors, these models face generalization and scalability bottlenecks.

Taken together, the two dark clouds: ineffective unified paradigm and the lack of strong prior knowledge, highlight the need for a new generation of unified models. In this work, we present **Muddit**, a MaskGIT-style unified discrete diffusion transformer equipped with a lightweight text decoder. By combining the strengths of parallel discrete diffusion and semantically rich visual priors from a pretrained Meissson text-to-image backbone (Bai et al., 2025), Muddit enables scalable, efficient, and flexible sampling while significantly improving alignment and quality across modalities and various tasks such as high-resolution text-to-image generation, image captioning, and visual question answering, as shown in Fig. 1 (c).

We systematically detail the training objective of unified discrete diffusion models, the masking strategy, and the shared inference sampling strategy across three tasks. Finally, we conduct comprehensive evaluations with current popular unified models on several benchmarks, including GenEval, CIDEr, VQAv2, GQA, MME, MMBench, and MMMU, demonstrating Muddit’s superior performance and efficiency, validating that the unexplored purely discrete diffusion approach can rival, or even surpass, much larger autoregressive-based unified models. While concurrent unified generation models (Yang et al., 2025b) often build upon a language modeling prior—leveraging pretrained dLLMs as the backbone—we instead take a visual-first approach. Muddit is built upon an image generation prior, offering a new path toward unifying vision and language tasks within a discrete diffusion framework. We hope that this work inspires a new trend for unified generative modeling, grounded in discrete diffusion, beyond the boundaries of traditional text-to-image generation (Bai et al., 2025).

## 2 METHOD

### 2.1 DISCRETE DIFFUSION WITH UNIFIED IMAGE AND TEXT PERSPECTIVE

In discrete diffusion, a sample  $x \in \mathcal{X}$  is treated as a one-hot vector  $\mathbf{x}$ , where  $\mathcal{X} = \{1, \dots, N\}$ . For language models,  $N$  equals the vocabulary size. While for image models,  $N$  is the number of discrete image token IDs obtained from a tokenizer or VQ codebook. At each diffusion step, we stochastically corrupt the tokens, gradually transforming the data distribution into a maximally entropic categorical prior; the generative model then learns to invert this corruption. Following recent works (Lou et al., 2023; Bai et al., 2025) that cast token corruption as a continuous-time Markov chain (CTMC) over the finite alphabet  $\mathcal{X}$ , we let

$$\frac{dp_t}{dt} = Q_t p_t, \quad (1)$$

where  $p_t \in \mathbb{R}^{N+1}$  is the distribution of  $x_t$ , the time-dependent matrix  $Q_t$  transports the data distribution  $p_0 \approx p_{\text{data}}$  to the maximally entropic “noise” distribution  $p_1 = p_{\text{stationary}}$ . We adopt the absorbing-state (masked) diffusion variant that has proved particularly effective in text modelling: every symbol can jump to a dedicated mask token  $\mathbf{m} = (\underbrace{0, \dots, 0}_N, 1)$  but never leaves it, i.e.  $\mathbf{m}$  is an

absorbing class.

**Forward posterior.** Marginalising  $\mathbf{x}$  gives

$$q(x_t | \mathbf{x}) = \text{Cat}(x_t | \alpha_t \mathbf{x} + (1 - \alpha_t) \mathbf{m}). \quad (2)$$

$\text{Cat}(\cdot)$  denotes a categorical distribution; it returns a one-hot token sampled from the probability vector inside the parentheses.  $\alpha_t \in [0, 1]$  is the *survival probability*, i.e. the probability that an individual token has not yet been masked by time  $t$ . Thus  $x_t$  equals the original clean token with probability  $\alpha_t$  and equals the mask token  $\mathbf{m}$  with probability  $1 - \alpha_t$ .

**Reverse process.** For any  $0 < s < t < 1$ , the CTMC induces an analytic posterior

$$q(x_s | x_t, \mathbf{x}) = \begin{cases} \text{Cat}(x_s | x_t), & x_t \neq \mathbf{m}, \\ \text{Cat}\left(x_s | \frac{(1 - \alpha_s)\mathbf{m} + (\alpha_s - \alpha_t)\mathbf{x}}{1 - \alpha_t}\right), & x_t = \mathbf{m}, \end{cases} \quad (3)$$

$x_t$  and  $x_s$  are the corrupted tokens at times  $t$  and  $s$  ( $s < t$ ). If  $x_t$  is already a real vocabulary token ( $x_t \neq \mathbf{m}$ ) it stays unchanged going backwards; otherwise, when  $x_t = \mathbf{m}$ , the distribution over  $x_s$  is a convex combination of the mask and the clean token  $\mathbf{x}$ , weighted by their respective survival probabilities  $\alpha_s$  and  $\alpha_t$ .

**Training Objective.** We employ a masked-token predictor  $x_\theta(x_t, \alpha_t) \approx \mathbf{x}$ , which leads to the continuous-time negative ELBO

$$\mathcal{L}_{\text{NELBO}} = \mathbb{E}_{q(x_t|\mathbf{x})} \left[ \int_0^1 \frac{\alpha'_t}{1 - \alpha_t} \log(x_\theta(x_t, \alpha_t) \cdot \mathbf{x}) dt \right], \quad (4)$$

where  $\alpha'_t = \frac{d\alpha_t}{dt}$  and  $\mathbf{x}$  is the one-hot vector of ground truth.  $x_\theta(x_t, \alpha_t) \in \mathbb{R}^{N+1}$  is the model’s predicted categorical probability vector for the clean token given the corrupted input  $(x_t, \alpha_t)$ ;  $\mathbf{x}$  is the one-hot ground-truth clean token.

During generation, we start from an all-mask sequence ( $t = 1$ ) and integrate the reverse CTMC towards  $t = 0$ , repeatedly replacing every masked position with the model’s categorical prediction. Because the corruption schedule and objective are *identical* for any discrete alphabet  $\mathcal{X}$ , the same diffusion backbone unifies text and image generation. In the following section, we present Muddit, a unified framework that leverages discrete diffusion to model the generation tasks for both text and image jointly.

## 2.2 MUDDIT

### 2.2.1 UNIFIED ARCHITECTURE

As shown in Fig. 2, our architecture comprises a text encoder  $E_{\text{txt}}$ , image encoder  $E_{\text{img}}$ , transformer generator  $G$ , sampler  $S$ , text decoder  $D_{\text{txt}}$ , and image decoder  $D_{\text{img}}$ . The generator  $G$  is a single MM-DiT model, following the dual-/single-stream design of FLUX (Labs, 2024). Importantly, the generator  $G$  is initialized from the Meissonic (Bai et al., 2025), which has been extensively trained for high-resolution text-to-image generation. This initialization brings in a strong pretrained image prior, capturing rich spatial structures and semantic correlations across image and text tokens, which significantly enhances sample quality and accelerates convergence in the multimodal setting. Consequently, the same MM-DiT predicts the masked tokens for both modalities, which produces a shared generator for text and image synthesis.

To reduce the computational cost of high-resolution imagery and lengthy captions, we quantize both modalities into a compact discrete space. A pre-trained VQ-VAE acts as the image encoder  $E_{\text{img}}$ , mapping pixels to codebook indices, while the CLIP text model, as  $E_{\text{txt}}$ , provides the text token embeddings. The MM-DiT predicts clean tokens in this shared space, which a lightweight linear head  $D_{\text{txt}}$  converts back to text tokens.

### 2.2.2 UNIFIED TRAINING

**Masking strategy.** We model the forward posterior in Eq. 2 of both modalities using time-dependent hyperparameters  $\alpha_t$ , with the mask ratio defined as  $\gamma_t = 1 - \alpha_t$ . While BERT (Devlin, 2018) employs a fixed mask ratio of 15%, this setting is suitable for token completion but insufficient for generation. To support generative tasks, the design of  $\gamma_t$  must satisfy the following criteria:

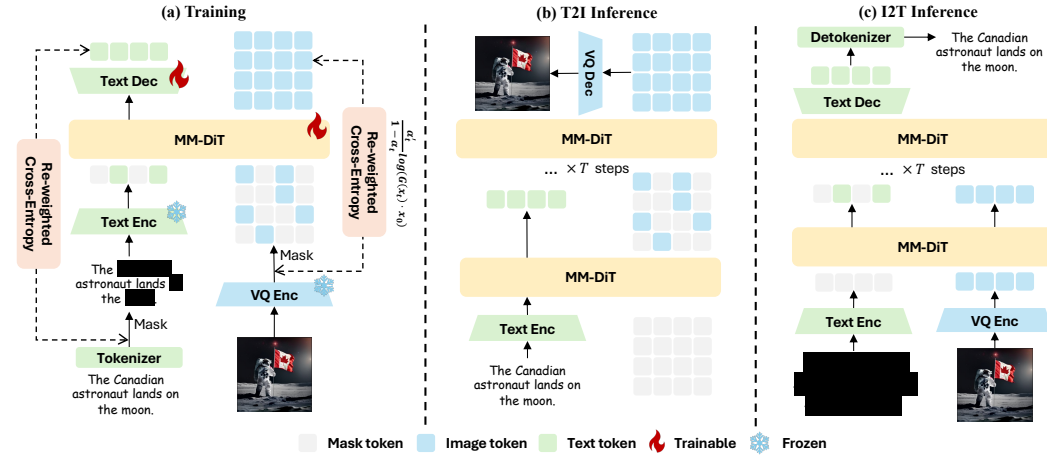


Figure 2: The training and inference architecture of Muddit. **(a)** During training, we randomly mask tokens from one of the two modalities. MM-DiT is trained to predict the masked tokens using a re-weighted cross-entropy loss, which jointly optimizes both the MM-DiT backbone and a lightweight text decoder. **(b)** In text-to-image inference, we initialize the image latent features using all-masked tokens and iteratively predict each latent token via MM-DiT. **(c)** In image-to-text inference, we similarly initialize all text tokens as masked and generate the text through the same iterative decoding process. Specifically for VQA tasks, we append mask token IDs to the end of the question and predict all masked token IDs as the final answer.

1.  $\gamma_t$  must be a continuous function, bounded between 0 and 1, for  $t \in [0, 1]$ .
2.  $\gamma_t$  should monotonically decrease with respect to  $t$ , with boundary conditions  $\gamma_0 \rightarrow 0$  (initially clean data) and  $\gamma_1 \rightarrow 1$  (masking all tokens).

Several strategies for masking and sampling have been proposed to meet these criteria (Chang et al., 2022). We adopt *cosine scheduling strategy*. During training, a timestep  $t \in [0, 1]$  is sampled from a truncated arccos distribution, with the density function:

$$\gamma_t = \frac{2}{\pi} (1 - (1 - t)^2)^{-\frac{1}{2}}. \quad (5)$$

During training, a mask ratio  $\gamma_t \in [0, 1]$  is randomly sampled for each modality  $\mathbf{x}_0$  (either image or text tokens), and the forward process (Eq. 2) is applied by randomly replacing clean tokens with mask tokens to obtain  $\mathbf{x}_t$ .

**Unified training objective.** Let  $\mathbf{c}$  denote the conditioning: the text embedding when synthesizing an image, or the image embedding when generating a caption. We randomly sample a mask ratio by Eq. 5. Then we corrupt the target sequence  $\mathbf{x}_0$  (image or text tokens) with the CTMC described in Eq. 1 and train a single masked-token predictor  $G(x_t, \alpha_t, \mathbf{c})$  to reconstruct  $\mathbf{x}_0$ . Both directions—text  $\rightarrow$  image and image  $\rightarrow$  text—share the identical continuous-time negative ELBO

$$\mathcal{L}_{\text{unified}} = \mathbb{E}_{q(x_t|\mathbf{x})} \left[ \int_0^1 \frac{\alpha'_t}{1 - \alpha_t} \log(G(x_t, \alpha_t, \mathbf{c}) \cdot \mathbf{x}) dt \right], \quad (6)$$

where all symbols are as in Eq. 4 but the  $G$  now receives the cross-modal condition  $\mathbf{c}$  as an additional input. **Key point:** switching from text  $\rightarrow$  image to image  $\rightarrow$  text merely changes the conditioning signal  $\mathbf{c}$ ; the loss Eq. 6 itself is unchanged. This symmetry keeps optimization identical across tasks and allows us to train a single parameter set jointly for both generation directions. During inference we again start from an all-mask sequence ( $t=1$ ) and integrate the reverse CTMC towards  $t=0$ , feeding in the desired condition  $\mathbf{c}$  to obtain either an image or a sentence from the same diffusion backbone.

### 2.2.3 UNIFIED INFERENCE

**Sampling strategy.** During inference, we apply the time-reversed posterior as defined in Eq. 3.

$$S(G, x_t, t) = p_\theta(x_s | x_t) = \begin{cases} \text{Cat}(x_s | x_t), & x_t \neq \mathbf{m}, \\ \text{Cat}\left(x_s | \frac{(1 - \alpha_s)\mathbf{m} + (\alpha_s - \alpha_t)G(x_t, \alpha_t, \mathbf{c})}{1 - \alpha_t}\right), & x_t = \mathbf{m}, \end{cases} \quad (7)$$

where  $\theta$  denotes the parameters of  $G$ ,  $\mathbf{c}$  is the multimodal condition, and  $\alpha_t$  in Eq. 5 is applied sequentially with  $t$  taking values  $1, \frac{T-1}{T}, \dots, \frac{1}{T}$ , where  $T$  is the total number of reverse steps. At each timestep  $t$ , Muddit predicts a fraction  $\gamma_{t+\frac{1}{T}} - \gamma_t$  of the masked tokens by  $G$  and update the masked tokens  $\mathbf{x}_t$  by  $S$ , continuing iteratively until all masked tokens are recovered. This dynamic approach offers several advantages over autoregressive methods, which require the model to learn conditional probabilities  $P(x_i | x_{<i})$  based on a fixed token ordering. In contrast, random masking with a variable ratio enables the model to learn  $P(x_i | x_\Lambda)$ , where  $\Lambda$  denotes an arbitrary subset of observed tokens. This flexibility is essential for parallel sampling, allowing multiple tokens to be predicted simultaneously rather than sequentially.

Our Muddit supports three tasks with a single generator  $G$  and sampler  $S$ : (i) text  $\rightarrow$  image, (ii) image  $\rightarrow$  text (captioning), and (iii) visual-question answering (VQA). The only change across tasks is the conditioning source  $\mathbf{c}$  provided to  $G$ ; the diffusion process and guidance logic are shared.

**(i) Text  $\rightarrow$  image.** Given a text prompt  $\mathbf{tp} \in \mathcal{T}$ , the text encoder  $E_{\text{txt}}$  produces a text token embedding  $\mathbf{c}_{\text{txt}} = E_{\text{txt}}(\mathbf{tp})$ . Starting from a fully masked sequence  $x_1$ , the generator produces logits

$$l_t = G(x_t, \alpha_t, \mathbf{c}_{\text{txt}}), \quad x_{t-\frac{1}{T}} = S(l_t, x_t, t), \quad (8)$$

for  $k = 1, \frac{T-1}{T}, \dots, \frac{1}{T}$ . After  $T$  steps we obtain visual tokens  $x_0$ , which the image decoder  $D_{\text{img}}$  converts to a pixel-space image  $I = D_{\text{img}}(x_0)$ .

**(ii) Image  $\rightarrow$  text.** For captioning, an input image  $I \in \mathcal{I}$  is tokenized by the image encoder  $E_{\text{img}}$ :  $\mathbf{c}_{\text{img}} = E_{\text{img}}(I)$ . The generator now conditions on the *visual* tokens while progressively decoding text:

$$l_t = G(x_t, \alpha_t, \mathbf{c}_{\text{img}}), \quad x_{t-\frac{1}{T}} = S(l_t, x_t, t), \quad (9)$$

yielding a text token sequence  $x_0$ , which  $D_{\text{txt}}$  maps to a caption  $\text{caption} = \text{Detokenize}(D_{\text{txt}}(x_0))$ .

**(iii) Image + question  $\rightarrow$  answer (VQA).** For visual-question answering we supply *both* an image and a question:  $\mathbf{c}_{\text{img}} = E_{\text{img}}(I)$  and  $\mathbf{c}_{\text{txt}} = E_{\text{txt}}(q)$ . They are concatenated and fed to the generator, which outputs logits over answer tokens  $x_k$ :

$$l_t = G(x_t, \alpha_t, [\mathbf{c}_{\text{img}}, \mathbf{c}_{\text{txt}}]), \quad x_{t-\frac{1}{T}} = S(l_t, x_t, t), \quad (10)$$

until the full answer  $a$  is produced and decoded by  $a = \text{Detokenize}(D_{\text{txt}}(x_0))$ .

**Classifier-free guidance.** At each decoding step, we apply the same guidance rule, independent of modality:

$$l_k \leftarrow G(z_k, \alpha_k, \mathbf{c}) + \lambda[G(z_k, \alpha_k, \mathbf{c}) - G(z_k, \alpha_k, \mathbf{c}_{\text{neg}})], \quad (11)$$

where  $z_k$  (image or text tokens) is the partial target sequence,  $\mathbf{c}$  is the *positive* condition (prompt, image, or image+question),  $\mathbf{c}_{\text{neg}}$  is the corresponding negative condition, and  $\lambda$  is the guidance scale. Because the loss, decoding schedule, and guidance operator are *identical* in all three scenarios—only the conditioning signal changes—our framework realises a genuinely unified multimodal generator.

## 3 EXPERIMENT

### 3.1 EXPERIMENTAL SETUP

**Implementation details.** We build Muddit on top of the open-sourced Meissson models (Bai et al., 2025). The MM-DiT backbone is initialized with pretrained weights, and a lightweight linear head



is added as a text decoder. Following Meissonic, we adopt the CLIP (Radford et al., 2021) as text encoder and VQ-VAE as image encoder and decoder, keeping them entirely frozen throughout all experiments. To support discrete denoising, we append a special `<mask>` token to CLIP’s vocabulary for text masking, while the image mask token is inherited directly from Meissonic’s initialization. We observe that, even without training, the `<mask>` embedding can already be predicted into a coherent sentence during training. Therefore, for simplicity, we freeze the `<mask>` embedding. During training, we use a constant learning rate of  $1 \times 10^{-4}$  and a weight decay of  $1 \times 10^{-2}$ . Gradient accumulation is applied in both pretraining and supervised fine-tuning, resulting in an effective batch size of 1024. We trained on 16 H100 GPUs for 5 days. During inference, we adopt the default Meissonic configuration, using cosine masking scheduling, 64 sampling steps, and a classifier-free guidance (CFG) scale of 9.0 and 1.5 for text-to-image and image-to-text generation, respectively.

**Training data.** We train Muddit in two stages using a combination of publicly available and internal datasets, including JourneyDB (Pan et al., 2023), LAION-Art (Schuhmann et al., 2022), CC12M (Changpinyo et al., 2021), and others. The final dataset is filtered based on aesthetic score, resolution, and aspect ratio, resulting in approximately 10 million image-text pairs. Both stages are optimized with the unified training objective defined in Eq. 6. Below, we describe the datasets and settings for each stage in detail.

1. **Pretraining.** We pretrain Muddit for 100K steps with a batch size of 1024, using the unified objective across both modalities. Text inputs are truncated to a maximum of 77 tokens, and images are resized to  $512 \times 512$ . The pretraining corpus consists of 8 million image-text pairs, re-captioned using Qwen2.5-VL-3B for improved consistency. Each batch is evenly split between text-to-image and image-to-text samples to enable joint training in both directions.
2. **Instruction tuning.** After pretraining, we fine-tune the model on a combination of 1 million instruction following datasets, including LLaVA-Instruct-150K, ALLaVA, SA-1B, and the VQAv2 training set. During this stage, only the answer portion of each prompt is masked. Additionally, we construct a curated dataset of 1 million high quality image-text pairs to support multi-task training on VQA and image generation. Following the task instructions embedded in each sample, Muddit learns to produce long-form answers, concise replies, and image captions via task-specific prompting.

We present both quantitative and qualitative results for the T2I and I2T tasks in the following sections. Additional experiments and ablation studies are provided in the Appendix.

### 3.2 TEXT-TO-IMAGE GENERATION

**Quantitative results.** Following prior work, we evaluate our  $512 \times 512$  model on GenEval (Ghosh et al., 2024) after supervised fine-tuning. Muddit attains an overall accuracy of 0.61, surpassing prior discrete diffusion models such as Monetico (0.44) and Meissonic (0.54), and closely matching Stable Diffusion 3 (0.62) with only 1B parameters. It further shows strong compositional reasoning (0.72 on “Two Objects”, 0.54 on “Counting”), and benefits from joint multimodal training, which enhances T2I performance. These results demonstrate the effectiveness of Muddit as the first unified discrete diffusion model for both text and image modalities.

**Qualitative results.** We present diverse generations from our model conditioned on rich textual prompts in Fig. 3. The outputs exhibit strong text-image alignment, capturing fine details in both realistic and imaginative scenes. Our model effectively renders complex structures, lighting, and textures across various domains.

### 3.3 IMAGE-TO-TEXT GENERATION

We present a comprehensive comparison of our model Muddit against other multimodal models across four benchmarks: MS-COCO (image captioning) (Lin et al., 2014), VQAv2 (Antol et al., 2015), MME (Fu et al., 2023), MMBench (Liu et al., 2024e), GQA (Hudson & Manning, 2019), and MMMU (Yue et al., 2024) in Tab. 2. Notably, Muddit is the first unified model to employ discrete diffusion for both text-to-image and image-to-text generation, demonstrating that this approach is not only viable but also highly competitive.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

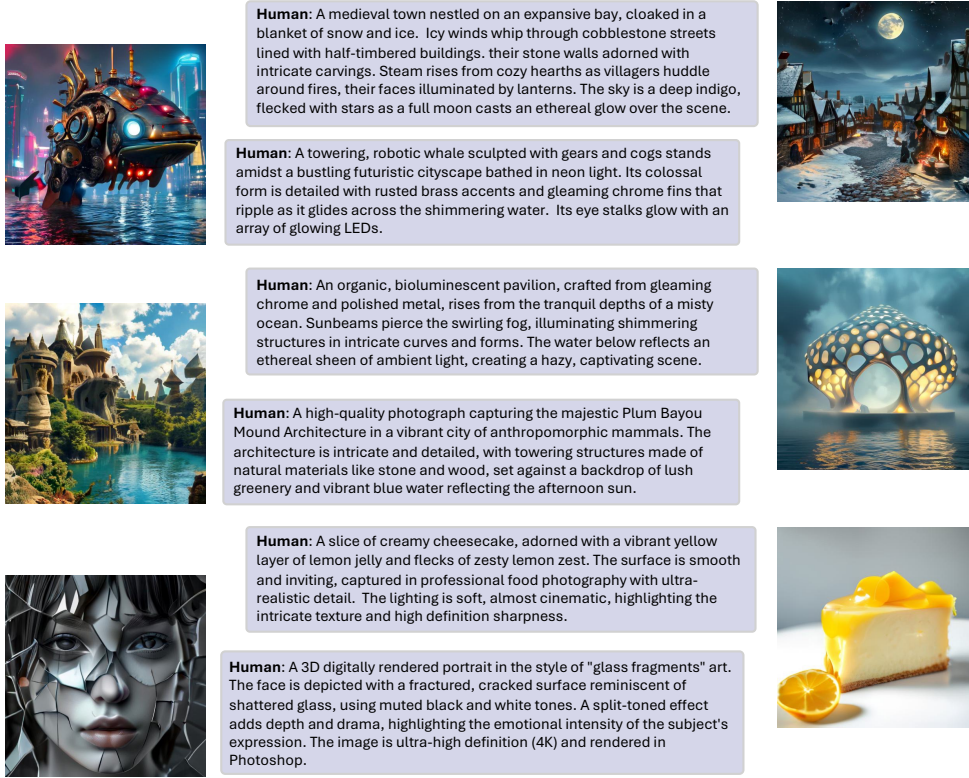


Figure 3: Samples of Text-to-Image Generation by Muddit.

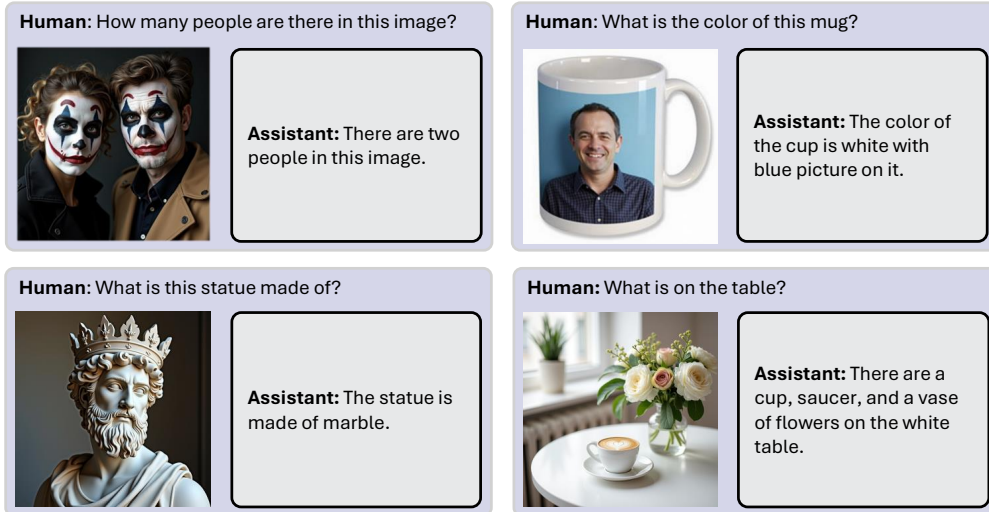


Figure 4: Samples of Visual Question Answering by Muddit.

**Quantitative comparison.** Muddit demonstrates strong performance across benchmarks despite its compact size (1B parameters). It achieves a CIDEr score of 59.7 on MS-COCO, surpassing diffusion-based baselines such as D-DiT (56.2) and Show-O (46.8–65.5). On VQAv2, it attains 67.7% accuracy, outperforming Show-O and D-DiT while approaching larger autoregressive models like LLaVA-Next (82.8%). Moreover, it reaches 1104.6 on MME, 28.4 on MMB, and 57.1 on GQA, underscoring its competitiveness across multimodal reasoning tasks. These results highlight the effectiveness of Muddit as a unified diffusion-based model that balances efficiency with high-quality task performance.

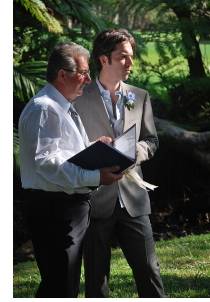
**Qualitative results.** We present example captions generated by our model across diverse scenarios in Fig. 5, including humans, animals, vehicles, and natural landscapes. The model demonstrates

Table 1: Evaluation of text-to-image generation performance on the GenEval (Ghosh et al., 2024).

Model	Text Gen Arch	Image Gen Arch	Params (B)	Overall	Objects ↑		Counting ↑	Colors ↑	Position ↑	Color ↑ Attribution
					Single	Two				
PixArt- $\alpha$ (Chen et al., 2024b)	-	Diffusion	0.6	0.48	0.98	0.50	0.44	0.80	0.08	0.07
SD 2.1 (Rombach et al., 2022)	-	Diffusion	0.9	0.50	0.98	0.51	0.44	0.85	0.07	0.17
DALL-E 2 (Ramesh et al., 2022)	-	Diffusion	6.5	0.52	0.94	0.66	0.49	0.77	0.10	0.19
SDXL (Podell et al., 2023)	-	Diffusion	2.6	0.55	0.98	0.74	0.39	0.85	0.15	0.23
DALL-E 3 (Betker et al., 2023)	-	Diffusion	-	0.67	0.96	0.87	0.47	0.83	0.43	0.45
SD 3 (Esser et al., 2024)	-	Diffusion	2	0.62	0.98	0.74	0.63	0.67	0.34	0.36
LWM (Liu et al., 2024c)	AR	AR	7	0.47	0.93	0.41	0.46	0.79	0.09	0.15
SEED-X (Ge et al., 2024)	AR	AR	17	0.49	0.97	0.58	0.26	0.80	0.19	0.14
Chameleon (Team, 2024)	AR	AR	7	0.39	-	-	-	-	-	-
Show-O (Xie et al., 2024)	AR	Discrete Diff.	1.3	0.68	0.98	0.80	0.66	0.84	0.31	0.50
Transfusion (Zhou et al., 2024)	AR	Diffusion	8	0.67	-	-	-	-	-	-
D-DiT (Li et al., 2024b)	Discrete Diff.	Diffusion	2	0.65	0.97	0.80	0.54	0.76	0.32	0.50
Monetico (512 × 512) (Bai et al., 2025)	-	Discrete Diff.	1	0.44	0.92	0.48	0.26	0.78	0.06	0.13
Meissonic (1024 × 1024) (Bai et al., 2025)	-	Discrete Diff.	1	0.54	0.99	0.66	0.42	0.86	0.10	0.22
UniDisc (512 × 512) (Swerdlow et al., 2025b)	Discrete Diff.	Discrete Diff.	1.4	0.42	0.92	0.47	0.15	0.67	0.13	0.19
Muddit (512 × 512)	Discrete Diff.	Discrete Diff.	1	0.61	0.98	0.72	0.54	0.82	0.19	0.41



**Assistant:** The image depicts a man standing in water, holding a black dog and smiling. The man is wearing a yellow and wet jacket, and he has a happy expression. The background consists of a calm body of water with possibly a lake.



**Assistant:** The image depicts two individuals standing in a lush, forested area. The person on the left is older and wearing a light blue shirt and dark pants, while the person on the right is dressed in a suit with a shirt and tie. Both appear to be engaged in a discussion or discussing something on the documents, with the person on the left holding a tablet.



**Assistant:** The image depicts a group of sheep walking along a road near a body of water, possibly a lake or a body of water, possibly a sea or sea. The road is lined with small rocks and a scattered vegetation, and a road is visible on the side of the road. In the background, there are rolling hills and a calm water with a fence.

**Assistant:** The image depicts a row of hot motorcycles parked in a parking lot. the motorcycle is a white model with a black and white color scheme, identifiable by its distinctive design and branding on the side. the motorcycle appears to be a modern model model, as evidenced by the branding and branding visible on the side. the background includes a residential building, some trees.



**Assistant:** The image depicts a brown horse standing amidst a lush, green forest. the horse appears to be a thoroughbred horse, characterized by its dense features and tail visible. it is equipped with a bridle and a decorative tag on its forehead. the background is filled with dense foliage and small plants, suggesting a natural, forested environment. the overall scene is serene.

Figure 5: Samples of Image-to-Text Generation by Muddit.

strong visual grounding and fine-grained descriptive ability, accurately capturing attributes such as clothing, expressions, background context, and object relationships. Fig. 4 illustrates our model’s ability to accurately answer visual questions across various domains, including object counting, color recognition, material identification, and compositional reasoning.

### 3.4 ABLATION STUDY AND ANALYSIS

**Analysis of the inference timesteps.** As shown in Tab. 5, performance generally improves with more diffusion steps, plateauing around  $T = 32$ . GenEval and CIDEr see large gains from  $T = 8$  to  $T = 32$ , with diminishing returns afterward. VQAv2 remains stable across timesteps, indicating that fewer steps suffice for discriminative tasks. Overall, a moderate number of steps provides a good balance between accuracy and efficiency.



Table 2: Evaluation of image captioning, visual question answering on multimodal benchmarks.

Model	Params (B)	Text Gen Arch	Image Gen Arch	MS-COCO CIDEr $\uparrow$	VQAv2 Acc. $\uparrow$	MME Acc. $\uparrow$	MMB Acc.	GQA Acc. $\uparrow$	MMMU Acc.
InternVL-2.0	8	AR	-	-	-	1648.1	81.7	61.0	49.3
LLaVA-Next	13	AR	-	-	82.8	1575.0	70.0	65.4	36.2
BLIP-2	13	AR	-	-	65.0	1293.8	-	41.0	34.4
QWEN-VL	7	AR	-	-	78.2	1487.5	-	57.5	35.9
OpenFlamingo	9	AR	-	65.5	43.5	-	-	-	28.7
Flamingo	9	AR	-	79.4	51.8	-	-	-	-
Chameleon	7	AR	AR	18.0	-	-	19.8	-	-
LWM	7	AR	AR	-	55.8	-	-	-	-
Show-O (256 $\times$ 256)	1.3	AR	Discrete Diff.	-	64.7	1014.9	-	54.2	-
Show-O (512 $\times$ 512)	1.3	AR	Discrete Diff.	-	69.4	1097.2	-	58.0	27.4
Transfusion	7	AR	Diffusion	29.0	-	-	-	-	-
D-DiT (256 $\times$ 256)	2	Discrete Diff.	Diffusion	-	59.5	897.5	-	55.1	-
D-DiT (512 $\times$ 512)	2	Discrete Diff.	Diffusion	56.2	60.1	1124.7	-	59.2	-
UniDisc	1.4	Discrete Diff.	Discrete Diff.	46.8	-	-	-	-	-
Muddit (512 $\times$ 512)	1	Discrete Diff.	Discrete Diff.	59.9	68.2	1107.4	28.4	57.5	27.6

Table 3: Impact of text loss weight. We apply the same text loss weight during both pretraining and instruction tuning.

Benchmark	0.2	0.4	0.6	0.8	1.0
GenEval	60.1	60.5	<b>61.6</b>	60.8	58.3
MS-COCO	51.4	52.1	<b>59.9</b>	58.8	59.4
VQAv2	62.7	66.2	68.2	68.4	<b>69.2</b>

**Analysis of the text loss weight.** As shown in Tab. 3, moderate text loss weights (around 0.6) yield the best overall performance. CIDEr and GenEval peak near this value, suggesting that both insufficient and excessive text weighting can harm generation quality. VQAv2 continues to improve with stronger text supervision but begins to plateau beyond 0.6. Overall, while discriminative tasks benefit from heavier textual guidance, generative tasks require a balanced mix of visual and textual signals—highlighting the importance of grounding language in multimodal learning.

**Analysis of joint training.** Joint optimization over both text-to-image (T2I) and image-to-text (I2T) objectives is essential. As shown in Tab. 4, joint training yields the highest GenEval score, outperforming both T2I-only and I2T-only variants. Notably, I2T-only causes GenEval to drop sharply from 61.6 to 28.3—more than a twofold decrease—while MS-COCO CIDEr remains nearly unchanged and VQAv2 declines only slightly. These results show that separating the objectives severely weakens cross-modal integration, underscoring the need for unified optimization to maintain strong multimodal coherence.

Table 4: Effect of joint training. We denote text-to-image as T2I and image-to-text as I2T, respectively.

Benchmark	T2I only	I2T only	Joint training
GenEval	59.3	28.3	<b>61.6</b>
MS-COCO	-	<b>60.1</b>	59.9
VQAv2	-	<b>69.1</b>	68.2

Table 5: Performance across different diffusion timesteps.

Sample steps	GenEval	CIDEr	VQAv2
T=8	51.6	43.6	53.9
T=16	58.5	59.3	57.4
T=24	59.3	59.4	62.3
T=32	<b>61.9</b>	59.7	65.4
T=40	61.7	60.1	66.8
T=64	61.1	59.9	<b>68.2</b>

### 3.5 THE SCALABILITY OF MUDDIT

To demonstrate the scalability of our approach, we curate roughly 10 million image-text pairs from LIAON-ART (Schuhmann et al., 2022), JourneyDB (Pan et al., 2023), and CC12M (Changpinyo et al., 2021). We filter out samples with an aesthetic score below 7, a height or width under 512 pixels, or an aspect ratio above 2. All images are re-captioned using Qwen2.5-VL 7B (Bai et al., 2023). We pretrain Muddit on this dataset with a batch size of 512 and a resolution of 1024, applying random masking to both image and text modalities. The image and text loss weights are set to 1.0 and 0.3, respectively. Training runs for 100K steps.

For instruction tuning, we collect about 6M samples from LLaVA-Instruct-150K (Liu et al., 2024d), ALLaVA LAION (Chen et al., 2024a), SA-1B (Kirillov et al., 2023), ART500K (Mao et al., 2017), ScienceQA (Lu et al., 2022), Chart2Text (Kantharaj et al., 2022), and VQAv2 (Antol et al., 2015). Muddit is then trained with a batch size of 512 at a resolution of 1024, with masking applied only to the answer text. We also add a 2M high-quality image dataset for high-quality fine-tuning. Further training configurations are provided in Tab. 6. All experiments are conducted on 16 H100 GPUs.

We evaluate the scaled Muddit model against other comparably sized unified models and state-of-the-art unified discrete diffusion models (Xin et al., 2025; Yang et al., 2025a), as shown in Tab. 7.

Table 6: Training hyperparameters across different training stages.

Hyperparameters	Stage-I (Pre-training)	Stage-II (Instruction-tuning)
Learning Rate	$1.0 \times 10^{-4}$	$1.0 \times 10^{-4}$
LR Scheduler	Constant	Constant
Weight Decay	0.01	0.01
Max Gradient Norm	10.0	10.0
Optimizer	AdamW ( $\beta_1 = 0.9$ , $\beta_2 = 0.999$ )	
Batch Size	512	512
Training Steps	100K	15K
Training GPUs	16×H100	16×H100
Gen. Resolution	1024	1024
Under. Resolution	1024	1024

Table 7: Quantitative comparison with other unified models.

Model	Params	Base model	Architecture	Data scale	Geneval w TTS	VQAv2	MME	MMMU
Lumina-DiMOO	8B	LLaDA	Discrete Diff.	80M	0.92	–	1534.2	58.6
MMaDA (512 × 512)	8B	LLaDA	Discrete Diff.	Unknown	0.66	76.7	1410.7	30.2
Show-O (512×512)	1.3B	Phi-1.5	AR + Discrete Diff	35M	–	69.4	1097.2	27.4
D-DiT (512×512)	2B	SD3-medium	Discrete Diff + Diff	40M	–	60.1	1124.7	–
Muddit (512×512)	1B	Meissonic	Discrete Diff	10M	0.64	68.2	1107.4	27.6
Muddit (1024×1024)	1B	Meissonic	Discrete Diff	16M	0.67	70.2	1139.2	28.7

Across established benchmarks, Muddit exhibits consistent improvements in both image generation and image understanding, empirically validating the scalability of our model. Furthermore, we compare Muddit with unified models of similar parameter sizes, all of which rely on hybrid architectures. Despite being trained on substantially less data, Muddit achieves superior performance.

We attribute this data efficiency to two key factors. First, our visual prior naturally maintains strong text-following capability for text-to-image generation, enabling robust alignment between image and text modalities. From the perspective of unified modeling, we prioritize cross-modal alignment over isolated single-modality ability, which allows Muddit to reach higher performance with less training data. Second, Muddit adopts a fully unified modeling paradigm: the model learns by predicting mask tokens based on context across all tasks (text-to-image and image-to-text). In contrast, hybrid architectures must simultaneously handle next-token prediction alongside velocity or mask prediction, and often introduce additional special tokens (*e.g.*,  $\langle \text{soi} \rangle$ ,  $\langle \text{eoi} \rangle$ ), which increases architectural complexity and hinders optimization.

## 4 CONCLUSION

In this work, we present Muddit, a unified generative framework that employs discrete diffusion to bridge text and image modalities. By unifying image and text generation within a single model, Muddit demonstrates strong performance across text-to-image, image-to-text, and VQA tasks. Notably, it outperforms or matches the capabilities of significantly larger autoregressive models, while enabling fast, parallel inference. Our results validate the effectiveness of discrete denoising as a general-purpose modeling strategy and highlight its potential to serve as a scalable backbone for future multimodal systems.

## REFERENCES

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pp. 2425–2433, 2015.
- Jinbin Bai, Tian Ye, Wei Chow, Enxin Song, Qing-Guo Chen, Xiangtai Li, Zhen Dong, Lei Zhu, and Shuicheng Yan. Meissonic: Revitalizing masked generative transformers for efficient high-resolution text-to-image synthesis. *ICLR*, 2025.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. URL <https://arxiv.org/abs/2308.12966>.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- 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.*, 2:3, 2023.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11315–11325, 2022.
- 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.
- Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model, 2024a.
- Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, Le Xue, Caiming Xiong, and Ran Xu. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset, 2025a. URL <https://arxiv.org/abs/2505.09568>.
- 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. In *ICLR*, 2024b.
- Liang Chen, Shuai Bai, Wenhao Chai, Weichu Xie, Haozhe Zhao, Leon Vinci, Junyang Lin, and Baobao Chang. Multimodal representation alignment for image generation: Text-image interleaved control is easier than you think. *arXiv preprint arXiv:2502.20172*, 2025b.
- Sixiang Chen, Jinbin Bai, Zhuoran Zhao, Tian Ye, Qingyu Shi, Donghao Zhou, Wenhao Chai, Xin Lin, Jianzong Wu, Chao Tang, et al. An empirical study of gpt-4o image generation capabilities. *arXiv preprint arXiv:2504.05979*, 2025c.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025d.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Runpei Dong, Chunrui Han, Yang 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.
- 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.

- Lijie Fan, Tianhong Li, Siyang Qin, Yuanzhen Li, Chen Sun, Michael Rubinstein, Deqing Sun, Kaiming He, and Yonglong Tian. Fluid: Scaling autoregressive text-to-image generative models with continuous tokens. *arXiv preprint arXiv:2410.13863*, 2024.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- 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.
- Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024.
- Jianyuan Guo, Hanting Chen, Chengcheng Wang, Kai Han, Chang Xu, and Yunhe Wang. Vision superalignment: Weak-to-strong generalization for vision foundation models. *arXiv preprint arXiv:2402.03749*, 2024.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pp. 16000–16009, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- 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*, pp. 6700–6709, 2019.
- Shankar Kantharaj, Rixie Tiffany Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. Chart-to-text: A large-scale benchmark for chart summarization. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4005–4023. Association for Computational Linguistics, 2022.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. Segment anything. In *ICCV*, 2023.
- Black Forest Labs. Announcing black forest labs, 2024. <https://blackforestlabs.ai/announcing-black-forest-labs/>.
- Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *Advances in Neural Information Processing Systems*, 37: 56424–56445, 2024a.
- Zijie Li, Henry Li, Yichun Shi, Amir Barati Farimani, Yuval Kluger, Linjie Yang, and Peng Wang. Dual diffusion for unified image generation and understanding, 2024b. URL <https://arxiv.org/abs/2501.00289>.
- Zijie Li, Henry Li, Yichun Shi, Amir Barati Farimani, Yuval Kluger, Linjie Yang, and Peng Wang. Dual diffusion for unified image generation and understanding. *arXiv preprint arXiv:2501.00289*, 2024c.
- 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 *ECV*, 2014.
- Dongyang Liu, Shitian Zhao, Le Zhuo, Weifeng Lin, Yu Qiao, Hongsheng Li, and Peng Gao. Lumina-mgpt: Illuminate flexible photorealistic text-to-image generation with multimodal generative pretraining, 2024a. URL <https://arxiv.org/abs/2408.02657>.

- Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024b.
- Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint*, 2024c.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 2024d.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *ECCV*, 2024e.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *NeurIPS*, 2022.
- Xinyin Ma, Runpeng Yu, Gongfan Fang, and Xinchao Wang. dkv-cache: The cache for diffusion language models. *arXiv preprint arXiv:2505.15781*, 2025.
- 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.
- Hui Mao, Ming Cheung, and James She. Deepart: Learning joint representations of visual arts. In *Proceedings of the 25th ACM international conference on Multimedia*, pp. 1183–1191. ACM, 2017.
- OpenAI. Addendum to gpt-4o system card: 4o image generation, 2025. URL <https://openai.com/index/gpt-4o-image-generation-system-card-addendum/>. Accessed: 2025-04-02.
- Junting Pan, Keqiang Sun, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Journeydb: A benchmark for generative image understanding, 2023.
- Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang Wang, Zhiyang Xu, Jiuhai Chen, Kunpeng Li, Felix Juefei-Xu, et al. Transfer between modalities with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025.
- 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.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- Christoph Schuhmann, Richard Vencu Beaumont, Romain Gordon Vencu, Clayton Coombes, Arun Katta, Robin Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next-generation image-text models. <https://laion.ai/blog/laion-5b/>, 2022. Accessed: 2025-09-25.



- Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative pretraining in multimodality. *arXiv preprint arXiv:2307.05222*, 2023.
- Alexander Swerdlow, Mihir Prabhudesai, Siddharth Gandhi, Deepak Pathak, and Katerina Fragkiadaki. Unified multimodal discrete diffusion. *arXiv preprint arXiv:2503.20853*, 2025a.
- Alexander Swerdlow, Mihir Prabhudesai, Siddharth Gandhi, Deepak Pathak, and Katerina Fragkiadaki. Unified multimodal discrete diffusion. *arXiv preprint arXiv:2503.20853*, 2025b. doi: 10.48550/arXiv.2503.20853.
- Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- Shengbang Tong, David Fan, Jiachen Zhu, Yunyang Xiong, Xinlei Chen, Koustuv Sinha, Michael Rabbat, Yann LeCun, Saining Xie, and Zhuang Liu. Metamorph: Multimodal understanding and generation via instruction tuning. *arXiv preprint arXiv:2412.14164*, 2024.
- 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.
- Chunwei Wang, Guansong Lu, Junwei Yang, Runhui Huang, Jianhua Han, Lu Hou, Wei Zhang, and Hang Xu. Illume: Illuminating your llms to see, draw, and self-enhance. *arXiv preprint arXiv:2412.06673*, 2024a.
- Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiyang Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024b.
- Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. Vila-u: a unified foundation model integrating visual understanding and generation. *arXiv preprint arXiv:2409.04429*, 2024.
- 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.
- Yi Xin, Qi Qin, Siqi Luo, Kaiwen Zhu, Juncheng Yan, Yan Tai, Jiayi Lei, Yuewen Cao, Keqi Wang, Yibin Wang, et al. Lumina-dimoo: An omni diffusion large language model for multi-modal generation and understanding. *arXiv preprint arXiv:2510.06308*, 2025.
- Ling Yang, Ye Tian, Bowen Li, Xincheng Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang. Mmada: Multimodal large diffusion language models. *arXiv preprint arXiv:2505.15809*, 2025a.
- Ling Yang, Ye Tian, Bowen Li, Xincheng Zhang, Ke Shen, Yunhai Tong, and Mengdi Wang. Mmada: Multimodal large diffusion language models. *arXiv preprint arXiv:2505.15809*, 2025b.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *CVPR*, 2024.
- Chuyang Zhao, Yuxing Song, Wenhao Wang, Haocheng Feng, Errui Ding, Yifan Sun, Xinyan Xiao, and Jingdong Wang. Monoformer: One transformer for both diffusion and autoregression. *arXiv preprint arXiv:2409.16280*, 2024.
- Mengyu Zheng, Yehui Tang, Zhiwei Hao, Kai Han, Yunhe Wang, and Chang Xu. Adapt without forgetting: Distill proximity from dual teachers in vision-language models. In *European Conference on Computer Vision*, pp. 109–125. Springer, 2024.
- 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.

# APPENDIX

## APPENDIX OVERVIEW

This appendix provides additional discussions, results, and analyses to complement the main paper. It is organized as follows:

- **Related Work** (Sec. A): We review unified multimodal models for understanding and generation, with a focus on autoregressive and diffusion-based paradigms, as well as recent advances in masked image modeling.
- **Additional Qualitative Results** (Sec. B): We present extended visualizations for several tasks, including image captioning, text-to-image generation, visual question answering, and image-guided text editing.
- **Additional Experimental Results** (Sec. C): We present more experimental results.
- **Additional Ablation Studies** (Sec. D): We present extended ablation studies.
- **Inference Time Analysis** (Sec. E): We analyze inference efficiency by comparing autoregressive decoding with discrete diffusion, providing FLOPs complexity and speed benchmarks.
- **Generated Results Step by Step** (Sec. F): We illustrate the reverse discrete diffusion process in detail, showing intermediate decoding steps and examples of progressive generation.
- **Discussion** (Sec. G): We reflect on the limitations of our approach and its broader impacts, including potential applications and risks of misuse.
- **Use of Large Language Models**: We clarify the role of large language models during paper preparation, emphasizing that they were only used for minor editing and polishing.

## A RELATED WORK

### A.1 UNIFIED MODELS FOR GENERATION AND UNDERSTANDING

The success of LLMs in language modeling has inspired efforts to extend unified generation to multimodal domains. However, the divergence between autoregressive and diffusion-based paradigms presents fundamental architectural trade-offs. Autoregressive models naturally handle language, and several works (Sun et al., 2023; Wang et al., 2024a; Tong et al., 2024; Ge et al., 2024; Dong et al., 2023; Chen et al., 2025b) extend this by connecting vision modules to LLMs via adapters or instruction tuning, with LLMs serving as planning modules that produce intermediate representations for image generation. While effective to some extent, these paradigms often exhibit limited interaction between text and image modalities and struggle with content consistency, particularly in image-to-image generation and complex instruction-based synthesis. To address these limitations, recent research explores unified generation models that integrate understanding and generation within a single architecture. We categorize these into four major paradigms (see Fig. 6):

**Fully Autoregressive:** Both text and image are tokenized into discrete sequences and modeled with an AR Transformer (Liu et al., 2024b; Team, 2024; Wu et al., 2024; Wang et al., 2024b; Chen et al., 2025d; Liu et al., 2024a; Guo et al., 2024; Zheng et al., 2024). These models achieve strong cross-modal generation but suffer from high latency due to sequential decoding.

**Text AR, Image Diffusion:** LLMs generate text tokens while image synthesis is delegated to pre-trained continuous diffusion backbones (Zhou et al., 2024; Zhao et al., 2024; Ma et al., 2024) or discrete diffusion (Xie et al., 2024). Though visually strong, these models are not truly unified, as they rely on separate architectures and token spaces.

**Image Diffusion, Text Discrete Diffusion:** Emerging models experiment with discrete diffusion for text and images (Li et al., 2024c), though many, like Dual-Diffusion, still use continuous diffusion for image synthesis, failing to realize true modality symmetry.

**Fully Discrete Diffusion:** Recent work like UniDisc (Swerdlow et al., 2025a) pioneers full-token discrete diffusion over shared Transformer backbones. These models support parallel sampling and native integration, but currently lag behind in generation fidelity and scale.

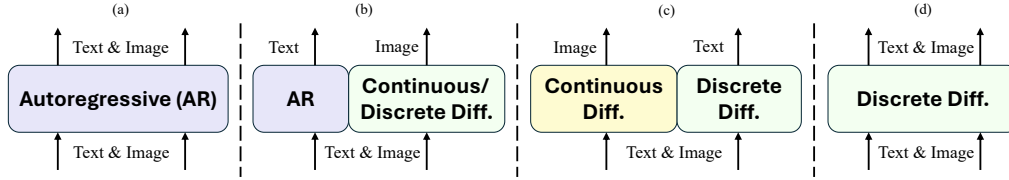


Figure 6: Four types of unified generative models. More details can be found in Sec. A.

Among these, the GPT-4o (OpenAI, 2025) model represents a significant advance as a unified multimodal generative system. However, its closed-source nature obscures critical architectural and training details, and its success may be largely attributable to scale rather than architectural novelty (Chen et al., 2025c).

## A.2 MASKED IMAGE MODELING

Masked Image Modeling (MIM) has emerged as a powerful self-supervised learning paradigm in computer vision, drawing inspiration from the success of Masked Language Modeling (MLM) in NLP, notably BERT (Devlin, 2018). The fundamental principle of MIM involves obscuring portions of an image, which could be raw pixels (MAE (He et al., 2022)), latent patches of pixels, or even discrete latent tokens (BEiT (Bao et al., 2021), MaskGIT (Chang et al., 2022)), and training a model, typically an autoencoder, to predict or reconstruct this missing information by leveraging the context provided by the visible parts.

MaskGIT (Chang et al., 2022) introduced parallel decoding via iterative token refinement, inspiring discrete diffusion models. Recent work such as RandomAR (Fan et al., 2024) and MAR (Li et al., 2024a) formalize this as random-order or masked autoregressive generation, blending AR and MIM principles. The major conceptual difference between RandomAR/MAR and MaskGIT is in the scanning order at inference time.

This class of techniques forms the conceptual foundation of discrete diffusion over tokenized spaces and plays a critical role in modern unified models. We will introduce discrete diffusion in the next section.

## A.3 RELATIONSHIP TO CONCURRENT WORK

Our main contribution is to show that a unified, *visual-prior* fully discrete diffusion model can be both effective and data-efficient on image understanding tasks, rather than only on text-to-image generation tasks. Regarding the distinction from discrete diffusion works (Swerdlow et al., 2025b; Yang et al., 2025a; Xin et al., 2025) we think that unified models should allow for multiple design choices. Our goal is to demonstrate that a visual-first, fully discrete diffusion backbone can be a practical and competitive alternative to the more common “LLM-first” unified paradigm, and we believe this is a fundamental design choice.

Concretely, prior unified discrete diffusion models such as UniDisc (Swerdlow et al., 2025b) are trained from scratch on multimodal data and therefore lack strong visual priors. As a result, they significantly underperform early diffusion baselines such as Stable Diffusion 1.5 (Rombach et al., 2022) and do not support vision–language reasoning tasks like VQA (Antol et al., 2015). In contrast, Mudit is the first unified discrete diffusion model built on top of a pretrained high-resolution text-to-image backbone (Meissonic), with a lightweight text decoder on top. This visual prior is not an implementation detail: it improves the scalability and generalization behavior of discrete diffusion through a text well-aligned visual backbone.

## B ADDITIONAL QUALITATIVE RESULTS

**Image-to-text Generation.** We present more examples for image-to-text generation in Fig. 7.

**Text-to-image Generation.** We present more examples for text-to-image generation in Fig. 8.

**Visual Question Answering.** We present more examples for visual question answering in Fig. 9. Muddit reliably identifies fine-grained attributes (*e.g.*, “blonde” hair), object categories (*e.g.*, “beagle”), and physical affordances (*e.g.*, answering “No” to crossing at a red light). Notably, it also handles commonsense reasoning and spatial localization, such as inferring traffic legality or locating vehicles on the street.

**Image-guided text editing.** Zero-shot text-guided image editing performance is already verified and presented in Meissonic (Bai et al., 2025). As the successor to Meissonic, we present Muddit’s performance on the image-guided text editing task, where the model completes a masked sentence based on the input image. As shown in Fig. 10, given a partially masked caption and an image, Muddit fills in the blanks with semantically and visually grounded phrases.

## C ADDITIONAL EXPERIMENTAL RESULTS

We provide a detailed breakdown of the MME benchmark results in the Tab. 8. Muddit demonstrates strong performance in existence, color, and scene understanding, while also exhibiting solid reasoning capabilities.

Table 8: Detailed MME results.

Category	Task	Score
Perception	Existence	135.00
	Count	78.33
	Position	53.33
	Color	140.00
	Posters	62.24
	Celebrity	56.18
	Scene	107.25
	Landmark	94.50
	Artwork	76.00
	OCR	52.50
	<b>Total</b>	<b>855.34</b>
Cognition	Commonsense Reasoning	78.57
	Numerical Calculation	90.00
	Text Translation	57.89
	Code Reasoning	57.50
	<b>Total</b>	<b>283.97</b>

## D ADDITIONAL ABLATION STUDIES

### D.1 ABLATION STUDY ON THE CFG FOR IMAGE-TO-TEXT GENERATION

As shown in Tab. 9. We report performance on MS-COCO captioning and VQAv2 benchmarks. Moderate CFG values (*e.g.*, 1.5) yield the best results, while higher scales lead to degraded performance.

## E INFERENCE TIME ANALYSIS

As shown in Fig. 13, autoregressive multimodal models are inherently limited by token-by-token decoding, which constrains their inference speed. Muddit overcomes this bottleneck with a parallel discrete diffusion decoder, reducing average latency to just 1.49 seconds, achieving a  $4\times$  to  $11\times$  speedup over competitive baselines ( $4.2\times$  faster than Qwen-2.5-VL,  $5.6\times$  than Show-o,  $8.1\times$  than BLIP-2, and  $10.9\times$  than LLaVA-1.6).

Besides, we present detailed FLOPs comparison between Autoregressive and Discrete Diffusion.

Table 9: Ablation study on the effect of classifier-free guidance (CFG) scale.

Dataset	CFG = 1	CFG = 1.5	CFG = 2	CFG = 2.5	CFG = 3
MS-COCO	57.2	59.9	58.2	51.3	47.2
VQAv2	65.8	68.2	64.7	55.4	49.2

Table 10: Comparison of model efficiency across different resolutions and steps. We report throughput for both text-to-image generation (images per second) and image-to-text tasks (tokens per second). Muddit achieves the best overall balance, matching the highest text-to-image throughput while significantly outperforming others in image-to-text speed.

Model	Image Res	Steps	Text-to-Image (img/s)	Image-to-Text (token/s)
Meissonic	1024	32	0.23	–
UniDisc	512	32	0.89	79.36
Monetico	512	32	1.00	–
D-DiT	512	28	0.62	26.89
Muddit	512	32	1.00	99.98

**Autoregressive (AR) without KV Cache:**

- At step  $t$ , the model attends over  $t$  previous tokens.
- Per-step attention FLOPs:  $O(t^2D)$ .
- Total FLOPs:

$$\sum_{t=1}^L O(t^2D) = O\left(D \sum_{t=1}^L t^2\right) = O\left(D \cdot \frac{L(L+1)(2L+1)}{6}\right) = O(L^3D)$$

**Autoregressive (AR) with KV Cache:**

- At step  $t$ , Q is computed for 1 token, and attends to  $t$  K/V keys.
- Per-step attention FLOPs:  $O(tD)$ .
- Total FLOPs:

$$\sum_{t=1}^L O(tD) = O\left(D \sum_{t=1}^L t\right) = O\left(D \cdot \frac{L(L+1)}{2}\right) = O(L^2D)$$

**Discrete Diffusion:**

- Each step updates the full sequence (length  $L$ ) in parallel.
- Per-step attention FLOPs:  $O(L^2D)$ .
- Total FLOPs:

$$T \cdot O(L^2D) = O(TL^2D), \quad T \ll L$$

While discrete diffusion may appear less efficient than autoregressive (AR) models with KV caching in terms of theoretical FLOPs, it offers a significant advantage over AR without caching—achieving an L/T speedup by updating the full token sequence in parallel over  $T$  iterations. In practice, the higher degree of parallelism leads to competitive, and often faster, inference speed compared to AR models, especially when considering real-world GPU throughput. As KV cache techniques for discrete diffusion are rapidly evolving (Ma et al., 2025), we expect further acceleration in the near future, narrowing the theoretical speed gap even with KV-cache AR baselines.

In Tab. 10, we compared Muddit against other non-autoregressive models, running all tests on a single A800 80 GB GPU. Muddit demonstrated a clear advantage in both image and text generation.



## F GENERATED RESULTS STEP BY STEP

Muddit frames text generation as reverse discrete diffusion over a fixed-length sequence of 77 token indices. At inference time, the model performs  $16 \leq T \leq 32$  denoising steps, starting from a maximally entropic prior where every token is masked. At each step  $t$ , a parameter-shared transformer  $G$  predicts a categorical distribution over all positions in parallel, and a sampler  $S$  selects the next sequence:

$$\mathbf{x}_{t-1} = S(G(\mathbf{x}_t, \mathbf{c}, t), \mathbf{x}_t, t), \quad t = T, \dots, 1, \quad (12)$$

where  $\mathbf{x}_t \in \mathbb{V}^{77}$  is the token sequence at step  $t$ , and  $\mathbf{c}$  denotes conditioning inputs. The logits can be tempered or top- $k$  filtered before sampling each token independently. The resulting sequence  $\mathbf{x}_{t-1}$  seeds the next step, enabling fast, parallel decoding without autoregressive constraints.

Because all positions are updated in parallel, Muddit preserves global syntactic and semantic structure throughout the reverse diffusion process—unlike left-to-right autoregressive models, which can only condition on past predictions. Empirically, as few as  $16 \leq T \leq 32$  steps are sufficient to approximate the natural language distribution with high fidelity. Thus, Muddit unifies diffusion generation with parallel decoding, effectively overcoming the serial bottleneck that limits conventional autoregressive multimodal models.

We present 2 examples in Fig. 11 and Fig. 12.

## G DISCUSSION

### G.1 LIMITATIONS

While Muddit advances discrete diffusion for unified multimodal generation, it still presents several limitations. First, due to its token-level discrete representation, the model may underperform continuous diffusion models in generating photorealistic or high-resolution images. Second, Muddit is initialized from a pretrained text-to-image foundation model, which offers strong visual priors but limits its capacity for rich text understanding and generation compared to the latest large language models. This makes it less suitable for tasks that require long-form understanding and generation or deep linguistic reasoning.

### G.2 BROADER IMPACTS

Muddit explores a new paradigm in multimodal generation by leveraging a strong visual prior as the backbone, in contrast to the prevailing trend of scaling large language models. This offers a complementary path toward efficient, grounded multimodal generation, particularly in vision-centric applications. The model’s ability to generate aligned visual and textual outputs in a fast, parallel manner could benefit downstream tasks, especially in completion-based scenarios such as masked captioning, image editing, and code implementation. However, as with all generative models, there remains a risk of misuse in synthetic content creation.

## H THE USE OF LARGE LANGUAGE MODELS

During the preparation of this paper, large language models were used only for language polishing and minor editing. All research ideas, methods, and experimental results were carried out entirely by the human authors.



**Assistant:** The image depicts a charming street scene with a cobblestone pavement and tall buildings. A white dog is resting. The dog is lying down near the bicycle on the street. To the left of the street, there is an old bicycle parked to the side of the street.

**Assistant:** The image is a black - and - white photograph featuring three individuals sitting together on a bench. They appear to be in an urban setting, as indicated by the trees and buildings in the background. The individuals are dressed in police uniforms, with one person looking behind the first individual from the second individual.



**Assistant:** The image depicts an elderly man standing on a city street, wearing a red jacket and a black cap talking on cell phone to his ear. The background shows a crowd of people walking around, some of whom are handbags and others visible in the distance. The crowd appears to be the middle of a busy area.

**Assistant:** The image depicts a close - up of a large, round pizza placed on a wooden table. The pizza is topped with cheese, and there are several pieces of basil leaves. In the background, there are other people sitting around.



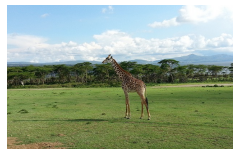
**Assistant:** The image depicts a tennis match on a blue court. the tennis match is set against a backdrop of a blue court. The player, wearing a white shirt and white shorts, is in mid action and holding a tennis racket, possibly preparing for a shot. the court has a row of seats.

**Assistant:** The image depicts a vibrant and colorful arrangement of umbrellas hanging from the branches of trees. The umbrellas come in color and shape, including pink, red, and yellow, blue, and purple. They are arranged in a circular pattern across the umbrellas, creating a visually appealing display. The background features a mix of green trees and foliage, suggesting an outdoor.



**Assistant:** The image depicts a bustling urban street scene with several notable elements. In the foreground, there is a tall green clock tower with a clock face and a pole, which stands prominently on the street. The clock has a clock face on it, indicating the time. In the background, there are several parked cars.

**Assistant:** The image depicts a person surfing on a surfboard on a large wave. The surfer is in mid-air, with his arms bent, navigating the balance of the wave. The background shows a clear blue sky, adding depth to the overall scene. The surfer seems to be in the middle of the the wave.



**Assistant:** The image depicts a giraffe standing on a grassy field in a rural outdoor setting. The background features a clear blue sky with some scattered clouds, and there are trees in the distant distance. The giraffe appears to be looking towards a distant horizon.

**Assistant:** The image depicts a stop sign in front of a grassy field with majestic mountains in the background. The sky is clear with a scattered clouds, and snow-capped mountains rise in the background. The scene appears to be in a rural setting.



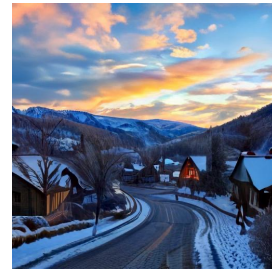
Figure 7: Image-to-text generated results.

**Human:** An editorial fashion photo portrait of a striking avant-garde model with bubble details. Iridescent art and pop surrealism influence the image's ethereal feel. Pastelaesthetic hues of soft pink, tangerine, and green create a dreamy backdrop. Seamless pink background, studio lighting emphasizes the model's silhouette against the soft pastel palette.



**Human:** Dark, heavy rainclouds gather over a rocky mountain range, obscuring the peaks. A weathered canvas tent flaps open in the wind, its silhouette stark against the gray backdrop. Water pours down in sheets, blurring the image and creating a misty effect.

**Human:** A quaint small town nestled amidst snow-capped hills, bathed in the soft, golden hues of dusk. The sky is a tapestry of twilight blue and orange, casting long shadows across cobblestone streets as frosted trees stand bare against the fading light. A sense of quiet peace emanates from the scene.



**Human:** An old cyanotype photograph capturing a serene lake at twilight. A majestic heron stands in the tranquil waters, its long legs poised, as the last golden rays of sunlight paint the sky with hues of blue and purple. The distant mountain range glows warmly with the soft, crepuscular light. Sharp focus on the heron and water reflection. Wide-angle lens captures a panoramic scene.

**Human:** A vibrant Japanese garden, inspired by Van Gogh's swirling brushstrokes. Vivid reds, yellows, blues, and greens dominate the scene. Cherry blossoms bloom in full splendor against a backdrop of ancient stone lanterns, rendered with lush foliage and dappled sunlight. The effect is heightened by 4K resolution and cinematic depth.



Figure 8: Text-to-image generation results.



Figure 9: Visual question answering results.

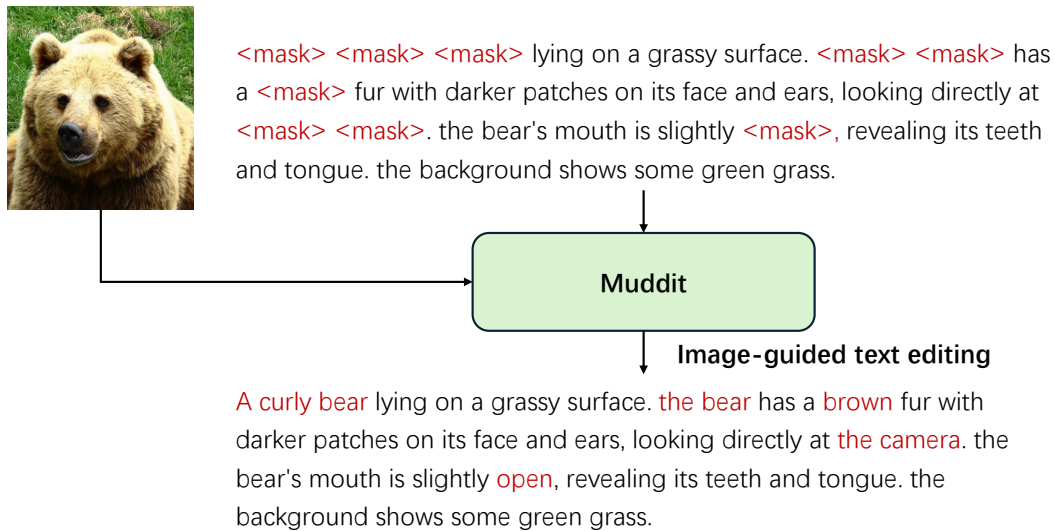


Figure 10: Image-guided text editing results.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241



**User:** Please describe this image.

**T=1**

<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>

<mask> <mask> depicts <mask> <mask> and whimsical <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask>, <mask> <mask> <mask> <mask> <mask> <mask>. <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>

The <mask> depicts a <mask> and whimsical <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask>, <mask> <mask> <mask> <mask> <mask> <mask>. the <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>,  
<mask> <mask> <mask> <mask> <mask> <mask>

The image depicts a surreal and whimsical scene in what appears to be a <mask> <mask>, possibly  
<mask> <mask> or a dining room. the floor is covered with <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask> <mask>  
<mask> <mask> <mask> <mask>, <mask> <mask> <mask> <mask> <mask>

The image depicts a surreal and whimsical scene in what appears to be a domestic setting, possibly a  
room or a dining room. the floor is covered with numerous pink flowers, adding a touch of <mask>  
<mask> <mask> <mask> <mask> <mask> <mask> <mask>, <mask> to the <mask> <mask>  
<mask> <mask> <mask> <mask> <mask>

The image depicts a surreal and whimsical scene in what appears to be a domestic setting, possibly a  
room or a dining room. the floor is covered with numerous pink flowers, adding a touch of <mask>. the  
petals are scattered throughout the room, adding to the dreamlike quality of the scene.

The image depicts a surreal and whimsical scene in what appears to be a domestic setting, possibly a  
room or a dining room. the floor is covered with numerous pink flowers, adding a touch of color. the  
petals are scattered throughout the room, adding to the dreamlike quality of the scene.

The image depicts a surreal and whimsical scene in what appears to be a domestic setting, possibly a  
room or a dining room. the floor is covered with numerous pink flowers, adding a touch of color. the  
petals are scattered throughout the room, adding to the dreamlike quality of the scene.

**T=0**

Figure 11: Image-to-text generated results in each step.





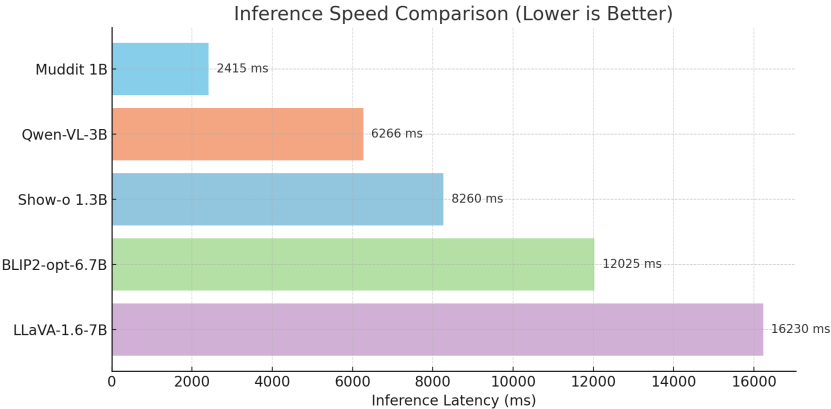


Figure 13: Inference speed comparison. We use 32 inference steps for Muddit and fix the sequence length to 77 across all models.