

# Transfer between Modalities with MetaQueries

Xichen Pan<sup>1,2</sup>, Satya Narayan Shukla<sup>1,†</sup>, Aashu Singh<sup>1</sup>, Zhuokai Zhao<sup>1</sup>, Shlok Kumar Mishra<sup>1</sup>, Jialiang Wang<sup>1</sup>, Zhiyang Xu<sup>1</sup>, Jiuhai Chen<sup>1</sup>, Kunpeng Li<sup>1</sup>, Felix Juefei-Xu<sup>1</sup>, Ji Hou<sup>1,†</sup>, Saining Xie<sup>2,†</sup>

<sup>1</sup>Meta, <sup>2</sup>New York University

<sup>†</sup>Equal advising

Unified multimodal models aim to integrate understanding (text output) and generation (pixel output), but aligning these different modalities within a single architecture often demands complex training recipes and careful data balancing. We introduce **MetaQueries**, a set of learnable queries that act as an efficient interface between autoregressive multimodal LLMs (MLLMs) and diffusion models. **MetaQueries** connects the MLLM’s latents to the diffusion decoder, enabling knowledge-augmented image generation by leveraging the MLLM’s deep understanding and reasoning capabilities. Our method simplifies training, requiring only paired image-caption data and standard diffusion objectives. Notably, this transfer is effective even when the MLLM backbone remains frozen, thereby preserving its state-of-the-art multimodal understanding capabilities while achieving strong generative performance. Additionally, our method is flexible and can be easily instruction-tuned for advanced applications such as image editing and subject-driven generation.

**Date:** April 9, 2025

**Correspondence:** [satyanshukla@meta.com](mailto:satyanshukla@meta.com), [jihou@meta.com](mailto:jihou@meta.com), [saining.xie@nyu.edu](mailto:saining.xie@nyu.edu)

**Project Page:** <https://xichenpan.com/metaquery>

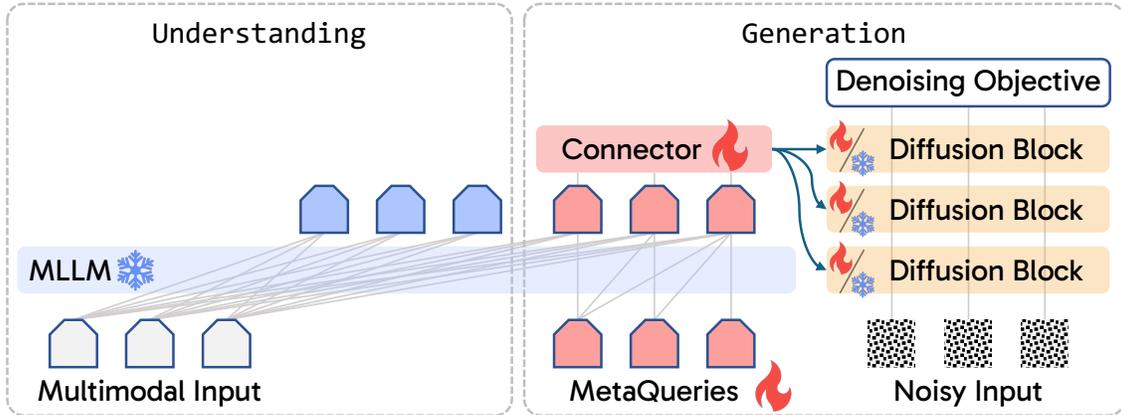


## 1 Introduction

The quest for unified multimodal models capable of both deep understanding (typically resulting in textual outputs) and rich generation (resulting in pixel outputs) holds immense promise. Such systems could unlock synergistic capabilities (OpenAI, 2025; Google, 2025), where understanding informs generation and vice versa. However, effectively connecting these different output modalities poses considerable challenges—*e.g.* how do we effectively transfer the latent world knowledge from the autoregressive multimodal LLM to the image generator? Although significant progress has been made, most published approaches (Ge et al., 2024; Sun et al., 2024b; Tong et al., 2024; Jin et al., 2024; Liu et al., 2024a; Team, 2024a; Xie et al., 2024; Wang et al., 2024; Wu et al., 2025a; Chen et al., 2025; Dong et al., 2024; Zhou et al., 2025; Shi et al., 2024) rely on carefully tuning base multimodal LLMs (MLLMs) to handle both understanding and generation tasks. This involves complex architectural design, data/loss balancing, multiple training stages, and other complex training recipes—without these, optimizing one capability could compromise the other.

In this paper, we aim to deliver the promise of unified models via a simpler philosophy: *Render unto diffusion what is generative, and unto LLMs what is understanding*. In other words, instead of building a monolithic system from scratch, we focus on effectively transferring capabilities between state-of-the-art, pre-trained models specialized for different output modalities. To operationalize this, we keep MLLMs frozen so they can focus on what they do best—understanding—while entrusting image generation to diffusion models. We then demonstrate that even under this frozen condition, the MLLM’s inherent world knowledge, strong reasoning, and in-context learning capabilities can indeed be transferred to image generation, provided the right architectural bridge is in place.

However, leveraging an MLLM—especially a frozen one—for both multimodal understanding and generation is far from straightforward. Although (frozen) LLMs have shown good performance as conditional text encoders in text-to-image generation (Zhuo et al., 2024; Xie et al., 2025; Ma et al., 2024), they are not compatible with many desired tasks in unified modeling, such as in-context learning or producing multimodal, interleaved output. The architectural bridge we design in this work is **MetaQuery** (Figure 1). **MetaQuery** feeds a set of



**Figure 1** Overview of our model. Blue tokens maintain SOTA multimodal understanding; MetaQueries are learnable queries that directly applied to frozen MLLMs to query out conditions for generation. The model is tuned using only denoising objective with paired data. The generative diffusion models can be either frozen or further instruction-tuned for advanced generation tasks.

learnable queries directly into a frozen MLLM to extract multimodal conditions for multimodal generation. Our experiments reveal that, even without fine-tuning or enabling bi-directional attention, the frozen LLM serves as a powerful feature resampler (Alayrac et al., 2022), producing high-quality conditions for multimodal generation. Training unified models with MetaQueries requires only a modest amount of paired image-caption data to connect these prompted conditions to any conditional diffusion model. Because the entire MLLM stays intact for understanding, the training objective remains the original denoising objective—just as efficient and stable as fine-tuning a diffusion model.

More specifically, previous unified models aim to train a single autoregressive transformer backbone to jointly model  $p(\text{text}, \text{pixels})$ . In contrast, we choose to use a  $\text{token} \rightarrow [\text{transformer}] \rightarrow [\text{diffusion}] \rightarrow \text{pixels}$  paradigm, which might share a high-level philosophy with the concurrent GPT-4o image generation system, as hinted at by OpenAI (2025). This approach composes the MLLM’s autoregressive prior with a powerful diffusion decoder, directly leveraging the frozen MLLM’s strong capability in modeling compressed semantic representations, thus avoiding the more challenging task of directly generating pixels.

To validate our approach, we conduct a series of controlled experiments, showing that MetaQuery<sup>1</sup> outperforms the use of a frozen MLLM purely as a conditional text encoder for image generation. Moreover, MetaQuery can match the performance of fully tuning the MLLM backbone, yet it is significantly more efficient. We also systematically investigate the training strategy, including the number of tokens and architectural configurations. With just 25M publicly available image-caption pairs, we are able to train a family of unified models that not only preserves state-of-the-art (SOTA) performance in image understanding, but also achieves SOTA-level results in text-to-image generation across multiple benchmarks.

The promise of unified modeling goes beyond handling multimodal understanding and text-to-image generation in parallel. A deeper synergy is expected—one that taps into advanced MLLM abilities like reasoning, internal knowledge, multimodal perception, and in-context learning to enhance generation. Our results show that our method draws on the frozen MLLM’s commonsense knowledge, achieving SOTA visual-commonsense generation on the CommonsenseT2I benchmark (Fu et al., 2024). Our approach also harnesses the built-in reasoning and in-context learning capabilities of frozen MLLMs, producing images from complex prompts—such as generating the United States flag in response to “The national flag of the country where Yellowstone National Park is located.” (See Figure 9 for examples.) We also benchmark this type of world knowledge reasoning capability on WISE (Niu et al., 2025) and demonstrate SOTA performance.

Finally, by connecting, preserving, and enhancing multimodal input with MetaQueries and a frozen MLLM backbone, our model can be further instruction-tuned for advanced generation tasks such as image editing and subject-driven generation. We show that this can be achieved both efficiently and effectively using a scalable

<sup>1</sup>For simplicity, we also use MetaQuery to represent our method.

data curation pipeline that directly leverages naturally occurring image pairs from web corpora, instead of depending on human-created pairs or synthetically generated data (Brooks et al., 2023; Hu et al., 2024a; Xiao et al., 2025). This natural supervision surprisingly unlocks several new capabilities beyond subject-driven generation, such as visual association and logo design (see Figure 8 for examples).

In summary, we explore a simple yet underexplored alternative to unified multimodal modeling. Our method, **MetaQuery**, bridges frozen MLLM backbones and diffusion models. Experiments show that this framework delivers all the capabilities once thought to require MLLM fine-tuning while being much easier to train. The main results and findings in this paper include:

- With **MetaQuery** and frozen MLLM backbones, we maintain SOTA multimodal understanding performance while enabling SOTA-level multimodal generation.
- **MetaQuery** can transfer the capabilities of MLLMs for reasoning- and knowledge-augmented image generation.
- **MetaQuery** can extract highly detailed visual conditions beyond semantic similarity from frozen MLLMs, enabling image reconstruction and editing tasks.
- Our method can be easily instruction-tuned even with a frozen MLLM backbone, enabling advanced multimodal generation tasks like subject-driven generation.

## 2 Related Work

*Unified understanding and generation models.* Next-token prediction has proven to be an effective approach for training models to understand language (Devlin, 2019; Brown et al., 2020) and multimodal content (Liu et al., 2024b). Recently, the community has witnessed numerous efforts to extend the success of multimodal understanding (Liu et al., 2024b) to multimodal generation by training LLM backbones to generate images at the same time. However, unlike adapting text-only LLMs (Touvron et al., 2023) to understand multimodal content with one single next text token prediction objective (Liu et al., 2024b), generating multimodal content requires a different set of training objectives. SEED-X (Ge et al., 2024), Emu (Sun et al., 2024b), and MetaMorph (Tong et al., 2024) learn to regress image features; LaVIT (Jin et al., 2024), LWM (Liu et al., 2024a), Chameleon (Team, 2024a), Show-o (Xie et al., 2024), EMU3 (Wang et al., 2024), and Janus (Wu et al., 2025a; Chen et al., 2025) auto-regressively predict next visual tokens; and DreamLLM (Dong et al., 2024), Transfusion (Zhou et al., 2025) employ diffusion objectives. However, these approaches necessitate tuning LLMs for generating both modalities, naturally posing challenges in multi-task balancing.

*Unified models with frozen LLMs.* Several studies have explored the use of frozen LLMs for multimodal understanding and generation. For instance, LMFusion (Shi et al., 2024) trains image generation expert feed-forward networks (FFNs) and query-key-value (QKV) modules in parallel with a frozen LLM backbone to deeply fuse input conditions and denoise visual outputs. However, this approach offers limited flexibility as it shares the same architecture as specific LLM backbones and requires training a separate set of generative modules for every single LLM backbone. This not only imposes more computational burden but also restricts the ability to leverage powerful pre-trained generative models. An earlier work, GILL (Koh et al., 2023), investigates feeding learnable tokens into frozen MLLMs. It employs a combined contrastive loss and regression loss for image retrieval and generation, rather than directly employing the denoising objective for more efficient training. Its application is restricted to contextual image generation and it does not systematically explore the impact of frozen MLLMs and learnable queries.

## 3 MetaQuery

In this work, we propose **MetaQuery**, which losslessly augments understanding-only MLLMs with multimodal generation capabilities while preserving their original architecture designs and parameters intact. We carefully analyze the impact of applying **MetaQuery** on image generation performance. Results show that a frozen MLLM can provide strong conditions for multimodal generation.

Methods	# of Tokens	MJHQ-30K FID ↓	GenEval ↑	DPG-Bench ↑
LLM last layer embedding*	-	7.49	0.55	78.41
Random queries	64	8.59	0.35	54.81
Learnable queries	64	7.43	0.56	75.35
Learnable queries	512	7.34	0.56	78.43

**Table 1** Study on different conditions for image generation. \* denotes the embeddings of input tokens. Learnable queries achieve comparable performance to using all hidden states and can even surpass them with more tokens.

Methods	Train LLM	Train DiT	MJHQ-30K FID ↓	GenEval ↑	DPG-Bench ↑
MLLM tuning	✓	✗	7.75	0.58	78.97
E2E tuning	✓	✓	6.28	0.61	79.39
Frozen MLLM	✗	✗	7.43	0.56	75.35
Frozen MLLM	✗	✓	6.06	0.61	76.66

**Table 2** Study on strategies for adapting MLLMs. The methods without training LLM do not suffer from multimodal understanding degradation. Frozen MLLM achieves comparable performance to full MLLM tuning, with slightly lower prompt alignment but slightly improved visual quality.

### 3.1 Architecture

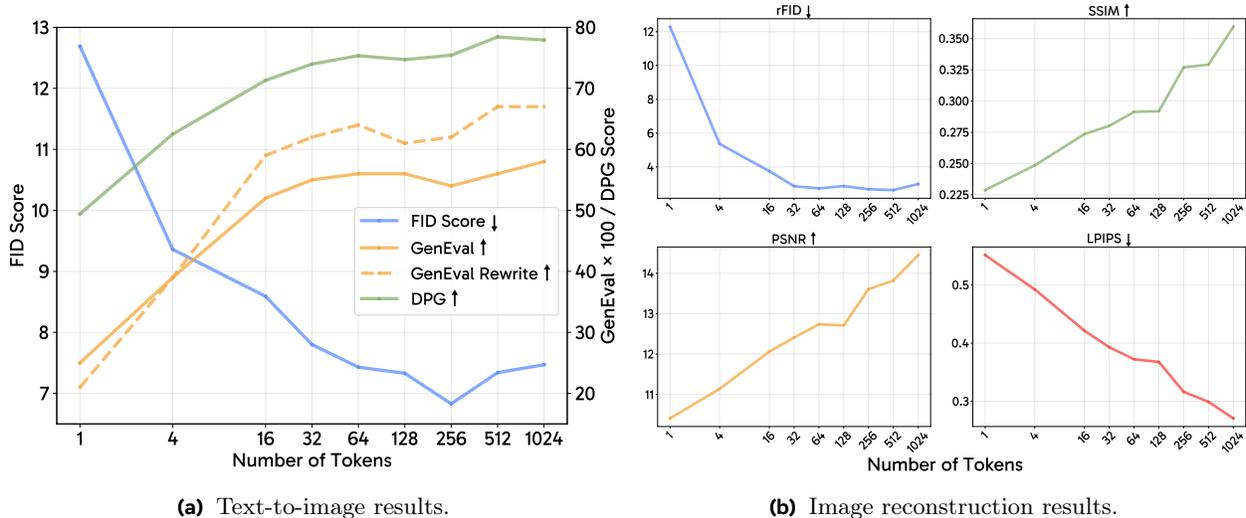
**MetaQuery** bridges frozen MLLMs with diffusion models. We use randomly initialized learnable queries  $\mathcal{Q} \in \mathbb{R}^{N \times D}$  to query out the conditions  $\mathcal{C}$  for generation.  $N$  is the number of queries and  $D$  is the dimension of the queries, which is the same as the MLLM hidden dimension. For simplicity and compatibility, we continue to use causal masking for the entire sequence rather than specifically enabling full attention for  $\mathcal{Q}$ . The conditions  $\mathcal{C}$  are then fed into a trainable connector to align with the input space of text-to-image diffusion models. These models can be arbitrary as long as they have a conditional input interface; we simply replace its original condition with our  $\mathcal{C}$ . The whole model is trained with the original generation objective on paired data. In this paper, we focus on image generation tasks, but the model can be easily extended to other modalities like audio, video, 3D, and more.

### 3.2 Design Choices

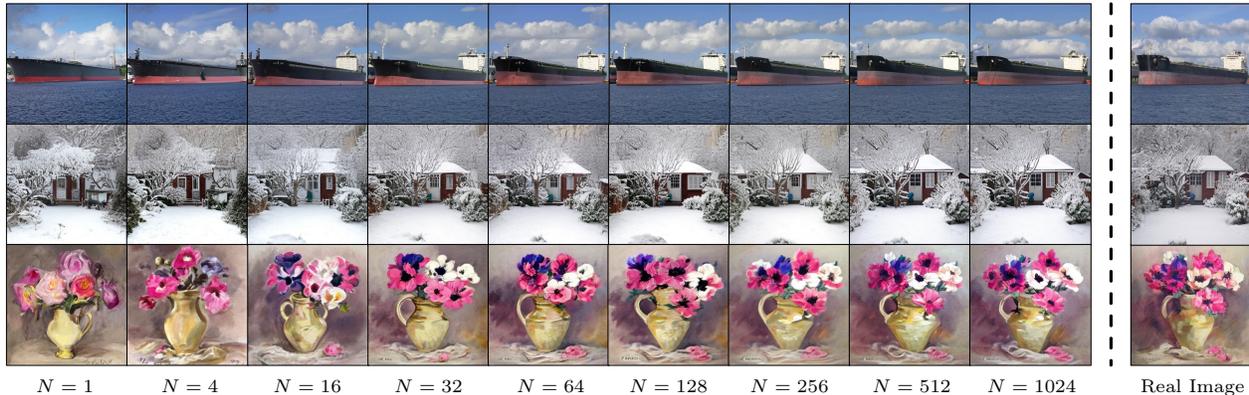
The proposed architecture involves two design choices: using **learnable queries** and keeping the **MLLM backbone frozen**. We explain the reasons why we adopted these choices and how they impact performance. For all experiments, unless otherwise specified, we use the same frozen LLaVA-OneVision-0.5B (Li et al., 2024a) MLLM backbone, frozen Sana-0.6B (Xie et al., 2025) diffusion model in 512 resolution, learnable queries with  $N = 64$  tokens, and a connector with a 24-layer transformer encoder. All models are trained on 25M publicly available image caption pairs for 4 epochs. We report FID score (Heusel et al., 2017) on MJHQ-30K (Li et al., 2024b) for visual aesthetic quality, and GenEval (Ghosh et al., 2023) and DPG-Bench (Hu et al., 2024b) (both without prompt rewriting) for prompt alignment, respectively.

*Learnable queries.* Many models like Lumina-Next (Zhuo et al., 2024), Sana (Xie et al., 2025), and Kosmos-G (Pan et al., 2024) use the (M)LLM’s last layer embedding of input tokens as image generation conditions. However, this approach is not ideal for unified models as it is not compatible with many desired tasks in unified modeling, such as in-context learning or producing multimodal, interleaved output (we provide more discussion and comparison with **MetaQuery** in Section 5.6). As shown in Table 1, using learnable queries with just  $N = 64$  tokens achieves image generation quality comparable to that of utilizing the last layer embedding of input tokens. While random queries produce acceptable FID scores, they struggle with prompt alignment, highlighting the importance of learnable queries. Additionally, since the last layer embedding setting naturally comes with a longer sequence length, we also tested learnable queries with  $N = 512$  tokens, which further improves performance and even outperforms the last layer embedding approach.

*Frozen MLLM.* Existing unified models train MLLMs to jointly model  $p(\text{text}, \text{pixels})$ , resulting in a more complicated training process and even downgraded understanding performance. **MetaQuery** keeps the original



**Figure 2** Study on the scaling of token numbers. As the number of tokens increases, text-to-image prompt alignment and image reconstruction results consistently improve.



**Figure 3** Visual samples for image reconstruction with different numbers of tokens.

MLLM architecture and parameters intact to preserve SOTA understanding capabilities. However, for multimodal generation, a key concern is whether *MetaQuery*'s performance with significantly fewer tunable parameters would be substantially worse than methods with full MLLM tuning. As shown in Table 2, frozen MLLMs achieve comparable performance to full MLLM tuning, with slightly lower prompt alignment but slightly improved visual quality. Tuning DiT can further improve performance for both settings. This suggests that *MetaQuery* is another possible training strategy, one that is simpler but also effective, as an alternative to fine-tuning the entire MLLM.

### 3.3 Training Recipe

Based on insights from our design choices, we further study key training options for the two main components of *MetaQuery*: learnable queries and connectors. This study examines the number of tokens and connector design. Unless otherwise specified, all experiments in this section use the same setup as described in Section 3.2.

*Number of tokens.* Many works (Wu et al., 2023; Pan et al., 2024; Ge et al., 2024) have employed learnable queries for condition extraction. However, they either set the number of tokens to match the fixed input sequence length of the image decoder (e.g.,  $N = 77$  for the CLIP (Radford et al., 2021) text encoder in Stable Diffusion v1.5 (Rombach et al., 2021)), or use an arbitrary fixed number like  $N = 64$  without further investigation. Given that modern diffusion models like Lumina-Next (Zhuo et al., 2024) and Sana (Xie

Architecture	# of Layers	Dims	# of Params	Rel. Wall Time	MJHQ-30K FID ↓	GenEval ↑	DPG-Bench ↑
Proj-Enc	6	2304	517M	1.06x	7.80	0.53	73.37
Proj-Enc	24	2304	2046M	1.23x	7.41	0.51	73.75
Enc-Proj	6	896	84M	1x	7.73	0.49	71.39
Enc-Proj	24	896	316M	1.06x	7.43	0.56	75.35

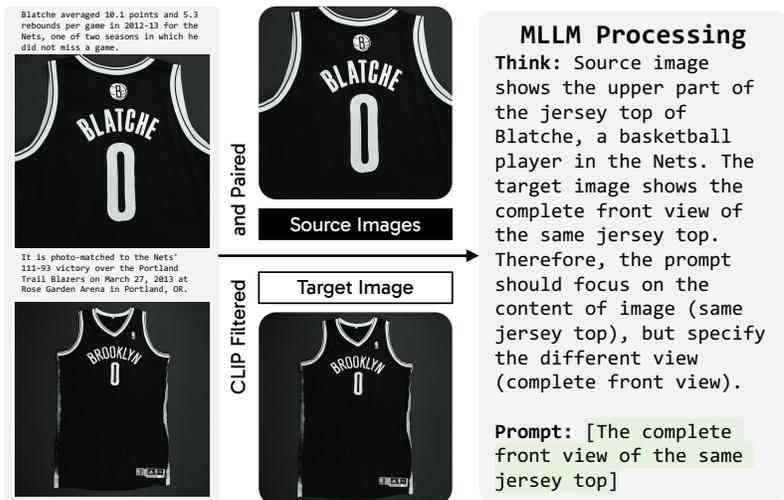
**Table 3** Study on connector design. Aligning the conditions first in the same dimension as the MLLM hidden states (Enc-Proj) is more effective and parameter-efficient.

et al., 2025) naturally accept variable-length conditions, determining the optimal number of tokens for learnable queries is crucial. In Figure 2, we provide a careful study of the number of tokens and observe promising scalability of *MetaQueries*. For text-to-image generation, visual quality begins to converge after 64 tokens, while more tokens consistently yield better prompt alignment. This is more evident for long captions, as GenEval with rewritten prompts increases more rapidly as the number of tokens increases. For image reconstruction, we observe that more tokens consistently improve the quality of reconstructed images (visual samples can be found in Figure 3). In our later experiments, we set the number of tokens to  $N = 256$  for all models, as it achieves a good balance between performance and efficiency.

*Connector design.* The connector is another important component in *MetaQuery*. We use the same architecture as the Qwen2.5 (Team, 2024b) LLM, but enable bi-directional attention for the connector. We study two different designs: Projection Before Encoder (Proj-Enc) and Projection After Encoder (Enc-Proj). Proj-Enc first projects the conditions into the input dimension of the diffusion decoder, then uses a transformer encoder to align the conditions. On the other hand, Enc-Proj first uses a transformer encoder to align the conditions in the same dimension as the MLLM hidden states, then projects the conditions into the input dimension of the diffusion decoder. As shown in Table 3, the Enc-Proj design achieves better performance than the Proj-Enc design while having fewer parameters.

## 4 Model Training

We train *MetaQuery* in two stages: the pre-training stage and the instruction tuning stage. Both training stages keep MLLMs frozen and fine-tune learnable queries, connectors, and diffusion models. We use three different MLLM backbones for different sizes: Base (LLaVA-OneVision 0.5B (Li et al., 2024a)), Large (Qwen2.5-VL 3B (Bai et al., 2025)), and X-Large (Qwen2.5-VL 7B (Bai et al., 2025)). We set the number of tokens to  $N = 256$  for all models, and utilize a 24-layer connector with Enc-Proj architecture. For image generation heads, we tested two different diffusion models: Stable Diffusion v1.5 (Rombach et al., 2021) and Sana-1.6B (Xie et al., 2025).



**Figure 4** Overview of instruction tuning data curation pipeline. We group images from web corpora based on caption similarity using the SigLIP (Zhai et al., 2023) model, then construct instruction-tuning data from these image pairs using an MLLM.

*Pre-training.* We pre-train our model on 25M publicly available image-caption pairs for 8 epochs with a learning rate of  $1e-4$  and a global batch size of 4096. The learning rate follows a cosine decay schedule with a 4,000-step warmup period before gradually decreasing to  $1e-5$ .

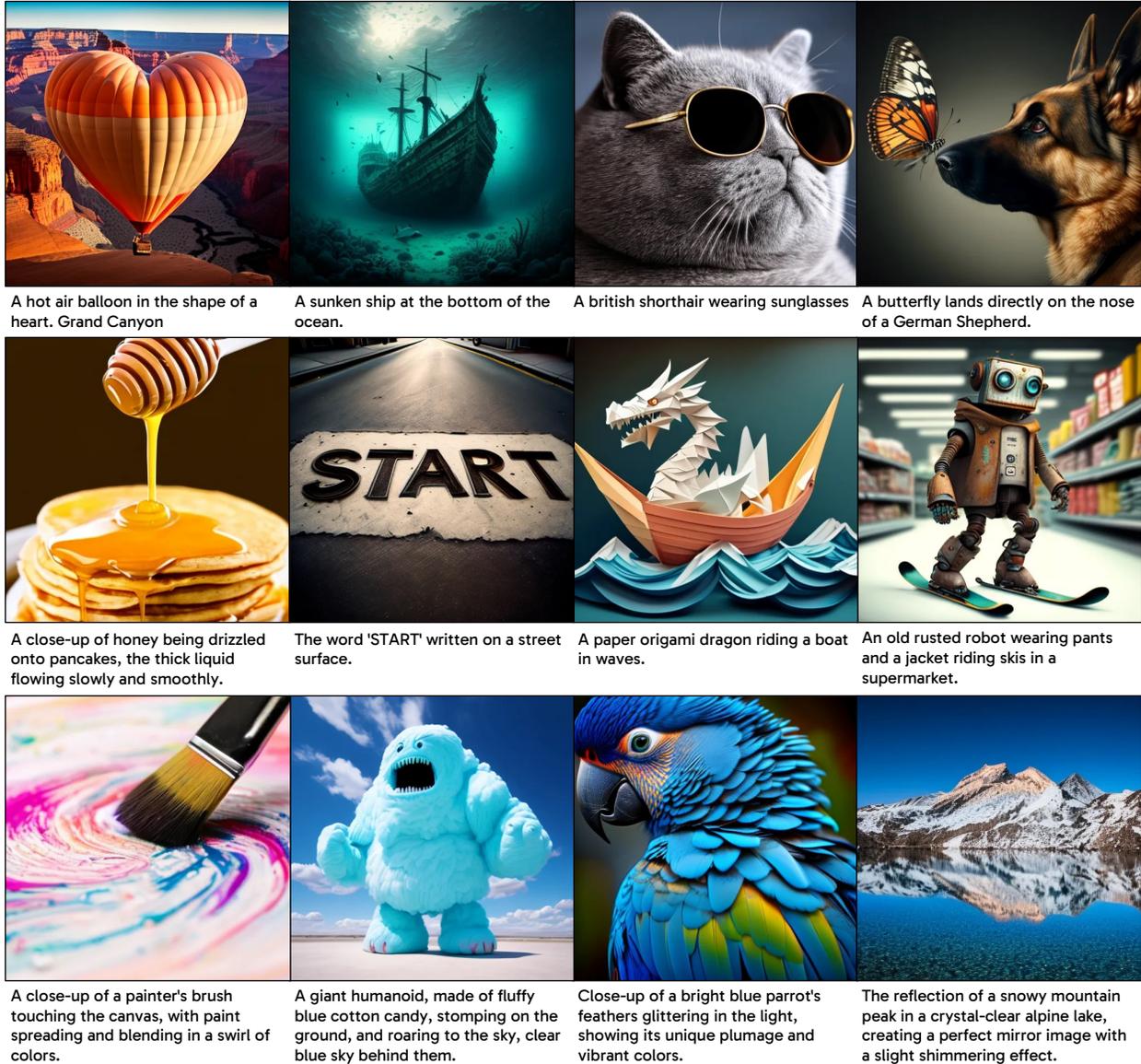
Methods	Base (M)LLM	MME-P	MMB	SEED	MMM	MM-Vet	COCO FID ↓	MJHQ FID ↓	GenEval ↑	DPG-Bench ↑
Emu	LLaMA 13B	-	-	-	-	-	11.66	-	-	-
DreamLLM	Vicuna 7B	-	-	-	-	36.6	8.46	-	-	-
Chameleon	From Scratch 7B	-	-	-	22.4	8.3	26.74	-	0.39	-
Show-o-512	Phi-1.5 1.3B	1097.2	-	-	26.7	-	9.24	15.18	0.68	-
VILA-U	LLaMA-2 7B	1401.8	-	59.0	-	33.5	-	7.69	-	-
Emu3	From Scratch 7B	-	58.5	68.2	31.6	37.2	12.80	-	0.66 <sup>†</sup>	80.60
MetaMorph	LLaMA-3 8B	-	75.2	71.8	-	-	11.8	-	-	-
TokenFlow-XL	Qwen-2.5 14B	1551.1	76.8	72.6	43.2	48.2	-	-	0.63 <sup>†</sup>	73.38
Transfusion	From Scratch 7B	-	-	-	-	-	8.70	-	0.63	-
LMFusion	LLaVA-Next 8B	1603.7	72.1	72.5	41.7	-	8.20	-	-	-
Janus	DeepSeek-LLM 1.5B	1338.0	69.4	63.7	30.5	34.3	8.53	10.10	0.61	-
JanusFlow	DeepSeek-LLM 1.5B	1333.1	74.9	70.5	29.3	30.9	-	9.51	0.63	80.09
Janus-Pro-1B	DeepSeek-LLM 1.5B	1444.0	75.5	68.3	36.3	39.8	-	14.33 <sup>‡</sup>	0.73	82.63
Janus-Pro-7B	DeepSeek-LLM 7B	1567.1	79.2	72.1	41.0	50.0	-	13.48 <sup>‡</sup>	0.80	84.19
MetaQuery-B	LLaVA-ov 0.5B	1238.0	58.5	66.6	31.4	29.1	8.91	6.28	0.74 <sup>†</sup>	80.04
MetaQuery-L	Qwen2.5-VL 3B	1574.3	78.6	73.8	53.1	63.2	8.87	6.35	0.78 <sup>†</sup>	81.10
MetaQuery-XL	Qwen2.5-VL 7B	1685.2	83.5	76.9	58.6	66.6	8.69	6.02	0.80 <sup>†</sup>	82.05

**Table 4** Quantitative results on multimodal understanding and generation benchmarks. We report the COCO FID with Stable Diffusion v1.5 (Rombach et al., 2021), and other metrics with Sana (Xie et al., 2025). <sup>†</sup> denotes rewritten prompts. <sup>‡</sup> denotes results tested by us under the same settings.

*Instruction tuning.* Furthermore, in this work, we rethink the data curation process for instruction tuning in image generation. All current methods rely on expert models to generate target images from source images and instructions (Ge et al., 2024; Xiao et al., 2025; Hu et al., 2024a). However, this approach is limited in scalability and may introduce biases, as the available expert models cover only a narrow range of image transformations. Inspired by MagicLens (Zhang et al., 2024), we construct instruction-tuning data using naturally occurring image pairs in web corpora. These corpora contain rich multimodal contexts with interleaved text and images on related subjects or topics. These image pairs often exhibit meaningful associations and specific relationships spanning a broad spectrum, from direct visual similarities to more subtle semantic connections (as shown in Figure 4). Such naturally occurring image pairs provide excellent and diverse supervision signals for instruction tuning. Based on this observation, we developed a data construction pipeline that mines image pairs and leverages MLLMs to generate open-ended instructions that capture their inter-image relationships. First, we collect grouped images from mmc4 (Zhu et al., 2023) core fewer-faces subset, where each image is accompanied by a caption. Using SigLIP (Zhai et al., 2023), we cluster images with similar captions (allowing up to 6 images per group, with a similarity threshold of 0.5). In each group, the image with minimum average similarity to the others is designated as the target, while the remaining images serve as source images. This process yields a total of 2.4M image pairs. Finally, we employ Qwen2.5-VL 3B (Bai et al., 2025) to generate instructions for each pair, describing how to transform the source images into the target image (See Appendix A for the detailed MLLM prompt). We experimented with instruction-tuning our Base size model on the proposed 2.4M dataset for 3 epochs, using the same learning rate schedule as in pre-training and a batch size of 2048.

## 5 Experiments

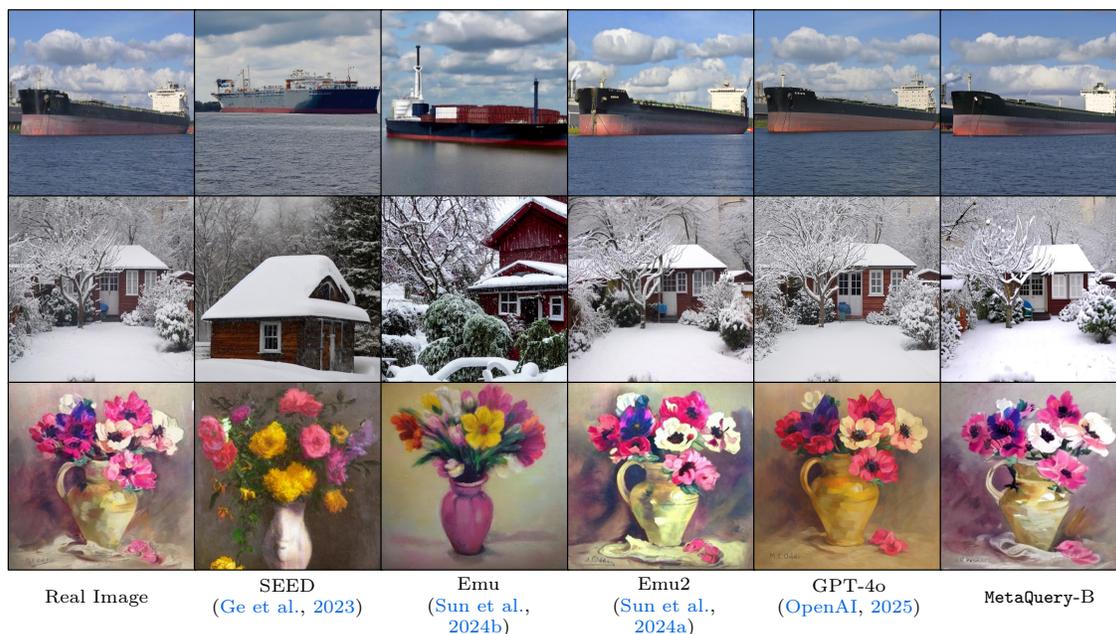
In this section, we first evaluate **MetaQuery** on various multimodal understanding and text-to-image generation benchmarks (Section 5.1). We demonstrate that **MetaQuery** can be trained to reconstruct input images (Section 5.2). This image reconstruction capability can be easily transferred to perform image editing (Section 5.3). Furthermore, we show that **MetaQuery** can be instruction-tuned to perform zero-shot subject-driven generation (Section 5.4). By leveraging our approach for collecting instruction tuning data from naturally existing image pairs, we also reveal that **MetaQuery** can unlock novel capabilities like visual association and logo design (also in Section 5.4). Additionally, we demonstrate that **MetaQuery** can benefit from the internal knowledge and reasoning capabilities of the frozen MLLM, overcoming common failures exhibited by other generation models (Section 5.5). Finally, we discuss the impact of different MLLM backbones and compare **MetaQuery**’s behavior with the baseline that uses MLLM last layer embeddings (Section 5.6).



**Figure 5** Qualitative results of MetaQuery. Prompts are from PartiPrompt (Yu et al., 2022), Sana (Xie et al., 2025) and Movie Gen Bench (Polyak et al., 2024).

## 5.1 Image Understanding and Generation

As shown in Table 4, our model family demonstrates strong capabilities across both understanding and generation tasks. Benefiting from the flexible training approach that allows us to leverage arbitrary SOTA frozen MLLMs, all of our models in different sizes exhibit competitive performance on all understanding benchmarks (Fu et al., 2023; Liu et al., 2023; Li et al., 2023a; Yue et al., 2024; Yu et al., 2023). In terms of image generation, MetaQuery achieves SOTA visual quality on MJHQ-30K (Li et al., 2024b). Given the fact that MetaQuery works with frozen MLLMs, we can naturally connect with an arbitrary number of diffusion models. Since the base Sana-1.6B (Xie et al., 2025) model is already fine-tuned on aesthetic data, we adopt Stable Diffusion v1.5 (Rombach et al., 2021) for COCO FID evaluation. Our results suggest that after adapting it to powerful MLLMs, we can achieve improved visual quality as indicated by the COCO FID score of 8.69. This also establishes a new SOTA COCO FID score among all Stable Diffusion v1.5-based unified models including MetaMorph (Tong et al., 2024) (11.8) and Emu (Sun et al., 2024b) (11.66).

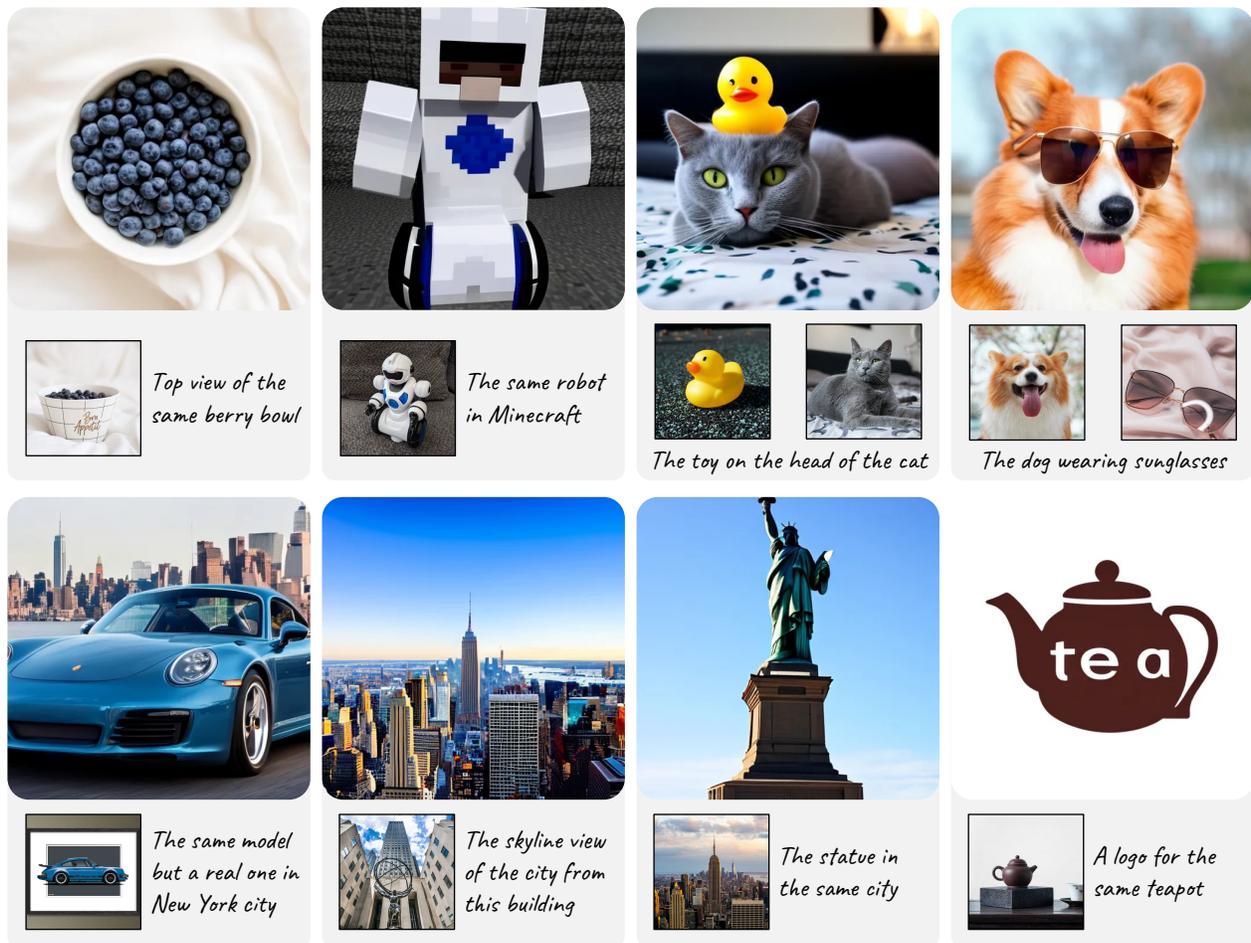


**Figure 6** Image reconstruction results. Results of SEED, Emu, and Emu2 are from Sun et al. (2024a).



**Figure 7** Image editing results. This capability can be easily transferred from image reconstruction after lightweight fine-tuning.

In terms of prompt alignment, **MetaQuery** also achieves competitive performance on GenEval (Ghosh et al., 2023) and DPG-Bench (Hu et al., 2024b), beating all diffusion model-based approaches including Transfusion (Zhou et al., 2025) and JanusFlow (Ma et al., 2025). We note that there is a performance gap between **MetaQuery** and Janus-Pro (Chen et al., 2025), which auto-regressively generates image tokens. We suggest that this gap may be due to the different failure modes of diffusion models and auto-regressive models: diffusion models usually fail to correctly follow the prompt, while auto-regressive models may suffer from more visual artifacts, which are difficult to quantify by GenEval and DPG-Bench. We tested the MJHQ-30K FID score of Janus-Pro under the same setting as ours and found that, in terms of visual quality and artifact control, **MetaQuery** is significantly better than Janus-Pro (see Appendix B for visual comparison). Additionally, we find that **MetaQuery** achieves much better world knowledge reasoning capability than Janus-Pro, which we will elaborate on in Section 5.5. We also found that when scaling up the size of frozen LLMs, the generation quality and prompt alignment also improves. **MetaQuery** provides a simple and principled way for leveraging the most advanced multimodal LLMs within a unified modeling framework. We also provide qualitative results in Figure 5 to illustrate the text-to-image generation capability of **MetaQuery**.



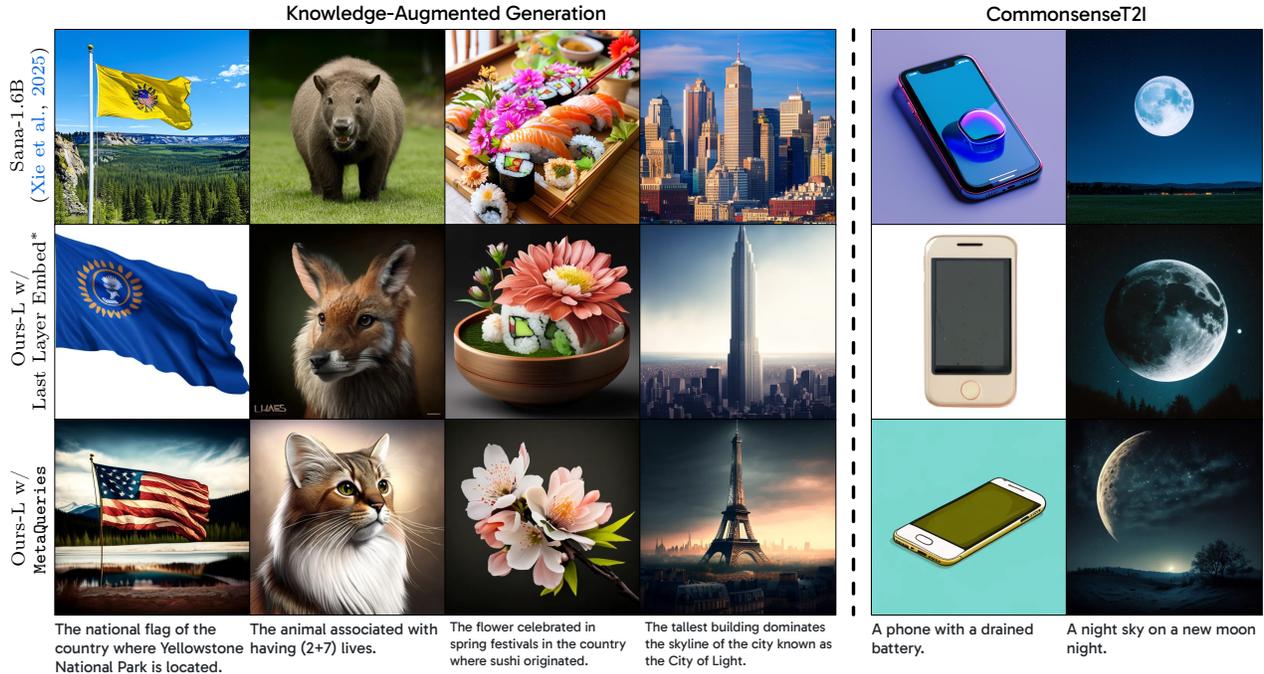
**Figure 8** Qualitative results for instruction tuning. Instruction-tuned MetaQuery achieves strong subject-driven capability (first row) and can even reason through the multimodal input to generate images (second row).

Methods	DINO Score $\uparrow$	CLIP-I Score $\uparrow$	CLIP-T Score $\uparrow$
Real Images (Oracle)	0.774	0.885	-
<i>fine-tuning</i>			
Textual Inversion (Gal et al., 2023)	0.569	0.780	0.255
DreamBooth (Ruiz et al., 2023)	0.668	0.803	0.305
BLIP-Diffusion (Li et al., 2023b)	0.670	0.805	0.302
<i>zero-shot &amp; test time tuning free</i>			
Re-Imagen (Chen et al., 2023)	0.600	0.740	0.270
BLIP-Diffusion (Li et al., 2023b)	0.594	0.779	0.300
Kosmos-G (Pan et al., 2024)	0.694	0.847	0.287
MetaQuery-B-Instruct	0.737	0.852	0.301

**Table 5** Subject-driven generation results on DreamBench (Ruiz et al., 2023).

## 5.2 Image Reconstruction

We demonstrate that MetaQuery can be easily fine-tuned for image reconstruction tasks with a frozen MLLM (See Appendix C for more details). As shown in Figure 6, we compare our fine-tuned MetaQuery-B with existing diffusion autoencoders from various unified models, which reconstruct images from predicted visual features. Since these unified models are not explicitly fine-tuned for image reconstruction, their results are directly decoded from the vision encoder’s output. Remarkably, even under this more constrained setup, our



**Figure 9** MetaQuery leverages frozen MLLMs for reasoning- and knowledge-augmented generation, overcoming the failure cases encountered in the base Sana model. \* denotes that the LLM last layer embeddings of input tokens are used for image generation; the model is in L size (Qwen2.5-VL 3B). This approach can be better than the base Sana model in some cases but fails to activate in-context learning to perform knowledge-augmented generation. Some of the test cases are from MetaMorph (Tong et al., 2024) and CommonsenseT2I (Fu et al., 2024).

fine-tuned MetaQuery-B can still achieve competitive performance, matching the best existing open-source model Emu2 (Sun et al., 2024a). When compared with GPT-4o (OpenAI, 2025), our model also achieves comparable quality.

### 5.3 Image Editing

As shown in Figure 7, we demonstrate that MetaQuery can transfer its image reconstruction capability to perform image editing. We keep the MLLM backbone frozen and fine-tune our pre-trained Base model for only 1,000 steps on publicly available image editing data. Qualitative results demonstrate that MetaQuery performs effectively in these image-editing scenarios.

### 5.4 Instruction Tuning

We show that after being instruction-tuned on the proposed 2.4M dataset in Section 4, MetaQuery can achieve impressive zero-shot subject-driven generation performance, producing coherent results even with multiple highly customized subjects (the first row of Figure 8). Using various supervision signals, the instruction-tuned MetaQuery-B model surprisingly unlocks novel capabilities like visual association and logo design that go beyond copy-pasting (the second row of Figure 8). For example, in the first case, the model identifies the specific model of the input Porsche 911 car image, then correctly generates a novel front view for that model. In the second case, the model recognizes the input image of Rockefeller Center and imagines the view of New York City from the top of the Rockefeller Center.

We also follow DreamBooth (Ruiz et al., 2023) by adopting DINO, CLIP-I, and CLIP-T scores to quantitatively evaluate our model on the DreamBench (Ruiz et al., 2023) dataset. As shown in Table 5, our MetaQuery-B-Instruct model achieves SOTA performance, outperforming existing models like Kosmos-G (Pan et al., 2024) that are explicitly trained on constructed substitution tasks for subject-driven generation.

Methods	Cultural	Time	Space	Biology	Physics	Chemistry	Overall
GPT-4o** (OpenAI, 2025)	<b>0.94</b>	<b>0.64</b>	<b>0.98</b>	<b>0.93</b>	<b>0.98</b>	<b>0.95</b>	<b>0.89</b>
<i>Text-to-Image Models</i>							
SD-v1-5 (Rombach et al., 2021)	0.34	0.35	0.32	0.28	0.29	0.21	0.32
SD-XL (Podell et al., 2023)	0.43	0.48	0.47	0.44	0.45	0.27	0.43
PixArt-Alpha (Chen et al., 2024)	0.45	0.50	0.48	0.49	0.56	0.34	0.47
playground-v2.5 (Li et al., 2024b)	0.49	0.58	0.55	0.43	0.48	0.33	0.49
SD-3.5-large (Esser et al., 2024)	0.44	0.50	0.58	0.44	0.52	0.31	0.46
FLUX.1-dev (Labs, 2024)	0.48	0.58	0.62	0.42	0.51	0.35	0.50
<i>Unified Models</i>							
show-o-512 (Xie et al., 2024)	0.28	0.40	0.48	0.30	0.46	0.30	0.35
vila-u-7b-256 (Wu et al., 2025b)	0.26	0.33	0.37	0.35	0.39	0.23	0.31
Emu3 (Wang et al., 2024)	0.34	0.45	0.48	0.41	0.45	0.27	0.39
Janus-1.3B (Wu et al., 2025a)	0.16	0.26	0.35	0.28	0.30	0.14	0.23
JanusFlow-1.3B (Ma et al., 2025)	0.13	0.26	0.28	0.20	0.19	0.11	0.18
Janus-Pro-1B (Chen et al., 2025)	0.20	0.28	0.45	0.24	0.32	0.16	0.26
Janus-Pro-7B (Chen et al., 2025)	0.30	0.37	0.49	0.36	0.42	0.26	0.35
MetaQuery-B	0.44	0.49	0.58	0.41	0.49	0.34	0.46
MetaQuery-L	0.56	0.57	0.62	0.48	0.63	0.42	0.55
MetaQuery-XL	0.56	0.55	0.62	0.49	0.63	0.41	0.55

**Table 6** Comparison of world knowledge reasoning on WISE (Niu et al., 2025). The test cases in WISE are similar to the knowledge-augmented generation ones in Figure 9. MetaQuery achieves SOTA performance and significantly outperforms all other unified models. \*\* Results are evaluated by Yan et al. (2025) on a random subset of 200 out of 1000 samples.

Methods	w/o Neg. Prompt	w/ Neg. Prompt
DALL-E 3 (Ramesh et al., 2021) w/ rewrite	40.17	N/A
SD-XL (Podell et al., 2023)	26.00	44.83
SD-3-medium (Esser et al., 2024)	26.17	47.17
FLUX.1-dev (Labs, 2024)	24.50	22.50
Sana-1.6B (Xie et al., 2025)	25.17	43.33
MetaQuery-B	27.33	51.50
MetaQuery-L	28.83	57.67

**Table 7** Comparison of visual commonsense reasoning capability on CommonsenseT2I (Fu et al., 2024).

## 5.5 Reasoning- and Knowledge-Augmented Generation

We show that the learnable queries can effectively leverage capabilities of the frozen LLM. This enables the model to better understand and follow complex prompts, including those requiring real-world knowledge and reasoning. As shown in Figure 9, for the left knowledge-augmented generation cases, MetaQuery-L can leverage world knowledge from the frozen MLLM and reason through the input question to generate the correct answer. For the right commonsense knowledge cases from CommonsenseT2I (Fu et al., 2024), the LLM provides better commonsense knowledge and enables MetaQuery to generate images that are consistent with the facts.

To quantitatively evaluate MetaQuery’s world knowledge reasoning capability, we employ the WISE (Niu et al., 2025) benchmark, which contains similar test cases to the knowledge-augmented generation examples shown in Figure 9. As demonstrated in Table 6, MetaQuery achieves SOTA performance, significantly outperforming all other unified models. Notably, before our work, existing unified models struggled to effectively leverage powerful MLLMs for reasoning and knowledge-augmented generation, resulting in inferior performance compared to text-to-image models. MetaQuery stands as the first unified model to successfully transfer the advanced capabilities of frozen MLLMs to image generation and exceed the performance of SOTA text-to-image models.

LLM Backbones	MJHQ-30K FID ↓	GenEval ↑	DPG-Bench ↑	CommonsenseT2I ↑
Qwen2.5-3B	6.20	0.79	81.34	56.00
Qwen2.5-3B-Instruct	6.36	0.79	81.12	54.33
Qwen2.5-VL-3B-Instruct	6.35	0.78	81.10	57.67

**Table 8** Comparison across different LLM backbones. Image generation capability is mostly orthogonal to multimodal understanding capability.

Methods	MJHQ-30K FID ↓	GenEval ↑	DPG-Bench ↑	WiScore ↑	CommonsenseT2I ↑
Ours-L w/ Last Layer Embed*	6.41	0.78	81.23	0.48	52.83
Ours-L w/ MetaQueries	6.35	0.78	81.10	0.55	57.67

**Table 9** Comparison between MetaQuery and LLM last layer embedding. \* denotes that the LLM last layer embeddings of input tokens are used for image generation. We observe comparable performance between MetaQuery and LLM last layer embedding on visual quality and prompt alignment. However, MetaQuery can activate in-context learning to perform knowledge-augmented generation, yielding much better performance on commonsense reasoning and world knowledge reasoning.

We also quantitatively evaluate MetaQuery’s commonsense reasoning capability on the CommonsenseT2I benchmark (Fu et al., 2024) in Table 7. For simplicity, we use CLIP (Radford et al., 2021) as the evaluator following their original implementation. Results show that MetaQuery significantly improves the performance of the base Sana model, achieving SOTA performance.

## 5.6 Discussion

*Comparison over different LLM backbones.* As shown in Table 8, to test the impact of employing different LLM backbones for MetaQuery, we carefully select a family of backbone models: pre-trained LLM (Qwen2.5-3B), instruction-tuned LLM (Qwen2.5-3B-Instruct), and instruction-tuned MLLM (Qwen2.5-VL-3B-Instruct). Both instruction-tuned models are initialized with the first pre-trained model checkpoint. Experimental results show that instruction tuning can achieve better (multimodal) understanding capabilities. However, the improvements are orthogonal to image generation performance when employed to provide multimodal generation conditions.

*Comparison with using last layer embeddings.* As shown in Table 1, our learnable queries approach achieves comparable image generation quality and prompt alignment to using the LLM’s last layer embeddings of input tokens. However, the last layer embedding method essentially treats the decoder-only LLM as a text encoder, which inherently limits its in-context learning capabilities. While this approach does improve upon the base Sana model in some cases as demonstrated in Figure 9, it struggles with the knowledge-augmented generation cases shown in the same figure. These cases require the LLM to first process and answer input questions before generating corresponding images, demanding in-context learning beyond what text encoders typically provide. This performance gap is quantitatively confirmed in Table 9, where MetaQuery significantly outperforms the last layer embedding approach on both WiScore and CommonsenseT2I benchmarks. Integrated natively with the LLM, MetaQuery naturally leverages its in-context learning capabilities, enabling the model to reason through questions and generate appropriate images.

## 6 Conclusion

We presented MetaQueries, a simple interface connecting MLLMs (for understanding) and diffusion decoders (for generation), effective even when the MLLM is frozen. This approach yields state-of-the-art understanding and generation performance with straightforward implementation. By enabling transfer between modalities, MetaQueries successfully channels MLLM knowledge and reasoning into multimodal generation. While effective, we hypothesize that bridging the remaining gap to leading proprietary systems may primarily involve further data scaling. We hope MetaQueries provides a powerful, accessible baseline for future unified multimodal model development.

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Siboz Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023.
- 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. In *NeurIPS*, 2020.
- 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*, 2024.
- Wenhu Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. In *ICLR*, 2023.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019.
- 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. In *ICLR*, 2024.
- 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 *ICML*, 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.
- Xingyu Fu, Muyu He, Yujie Lu, William Yang Wang, and Dan Roth. Commonsense-t2i challenge: Can text-to-image generation models understand commonsense? In *COLM*, 2024.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023.
- 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.
- 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. In *NeurIPS*, 2023.
- Google. Experiment with gemini 2.0 flash native image generation, 2025. <https://developers.googleblog.com/en/experiment-with-gemini-20-flash-native-image-generation/>.
- 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. In *NeurIPS*, 2017.
- Hexiang Hu, Kelvin CK Chan, Yu-Chuan Su, Wenhu Chen, Yandong Li, Kihyuk Sohn, Yang Zhao, Xue Ben, Boqing Gong, William Cohen, et al. Instruct-imagen: Image generation with multi-modal instruction. In *CVPR*, 2024a.
- Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024b.

- Yang Jin, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Bin Chen, Chenyi Lei, An Liu, Chengru Song, Xiaoqiang Lei, et al. Unified language-vision pretraining with dynamic discrete visual tokenization. In *ICLR*, 2024.
- Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. Generating images with multimodal language models. In *NeurIPS*, 2023.
- Black Forest Labs. Flux.1, 2024.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024a.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023a.
- Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2. 5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv preprint arXiv:2402.17245*, 2024b.
- Dongxu Li, Junnan Li, and Steven CH Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. In *NeurIPS*, 2023b.
- Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2024b.
- 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? *arXiv preprint arXiv:2307.06281*, 2023.
- Bingqi Ma, Zhuofan Zong, Guanglu Song, Hongsheng Li, and Yu Liu. Exploring the role of large language models in prompt encoding for diffusion models. In *NeurIPS*, 2024.
- 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. In *CVPR*, 2025.
- Yuwei Niu, Munan Ning, Mengren Zheng, Bin Lin, Peng Jin, Jiaqi Liao, Kunpeng Ning, Bin Zhu, and Li Yuan. Wise: A world knowledge-informed semantic evaluation for text-to-image generation. *arXiv preprint arXiv:2503.07265*, 2025.
- OpenAI. Introducing 4o image generation, 2025. <https://openai.com/index/introducing-4o-image-generation/>.
- Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhui Chen, and Furu Wei. Kosmos-g: Generating images in context with multimodal large language models. In *ICLR*, 2024.
- 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.
- Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2021.
- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023.
- Weijia Shi, Xiaochuang Han, Chungting Zhou, Weixin Liang, Xi Victoria Lin, Luke Zettlemoyer, and Lili Yu. Llamafusion: Adapting pretrained language models for multimodal generation. *arXiv preprint arXiv:2412.15188*, 2024.

Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyong Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *CVPR*, 2024a.

Quan Sun, Qiyong Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative pretraining in multimodality. In *ICLR*, 2024b.

Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024a.

Qwen Team. Qwen2.5: A party of foundation models, 2024b.

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.

Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiyong Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024.

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. In *CVPR*, 2025a.

Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. *arXiv preprint arXiv:2309.05519*, 2023.

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. In *ICLR*, 2025b.

Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. In *CVPR*, 2025.

Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion transformers. In *ICLR*, 2025.

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.

Zhiyuan Yan, Junyan Ye, Weijia Li, Zilong Huang, Shenghai Yuan, Xiangyang He, Kaiqing Lin, Jun He, Conghui He, and Li Yuan. Gpt-imgeval: A comprehensive benchmark for diagnosing gpt4o in image generation. *arXiv preprint arXiv:2504.02782*, 2025.

Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. In *TMLR*, 2022.

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.

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 *CVPR*, 2024.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023.

Kai Zhang, Yi Luan, Hexiang Hu, Kenton Lee, Siyuan Qiao, Wenhui Chen, Yu Su, and Ming-Wei Chang. Magiclens: Self-supervised image retrieval with open-ended instructions. In *ICML*, 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. In *ICLR*, 2025.

Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal C4: An open, billion-scale corpus of images interleaved with text. In *NeurIPS*, 2023.

Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenze Liu, Lirui Zhao, Fu-Yun Wang, Zhanyu Ma, et al. Lumina-next: Making lumina-t2x stronger and faster with next-dit. *arXiv preprint arXiv:2406.18583*, 2024.

# Appendix

## A Data Curation Details

For the data curation part, we use `Qwen/Qwen2-VL-7B-Instruct`<sup>2</sup> as our MLLM, The system prompt we are using is:

Based on the provided of one or multiple source images, one target image, and their captions, create an interesting text prompt that can be used with the source images to generate the target image. This prompt should include:

- one general and unspecific similarity shared with the source images (same jersey top, similar axe, similar building, etc).
- all differences that only the target image has.

This prompt should NOT include:

- any specific details that would allow generating the target image independently without referencing the source images.

Remember the prompt should be concise and short. The generation has to be done by combining the source images and text prompts.

## B Qualitative Comparison with SOTA Open-Source Model on Text-to-Image Generation

We provide a qualitative comparison with Janus-Pro-7B (Chen et al., 2025) on MJHQ-30K (Li et al., 2024b) in Figure 10. We can see that MetaQuery-XL follows the prompt better and generates more visually appealing images than Janus-Pro-7B.

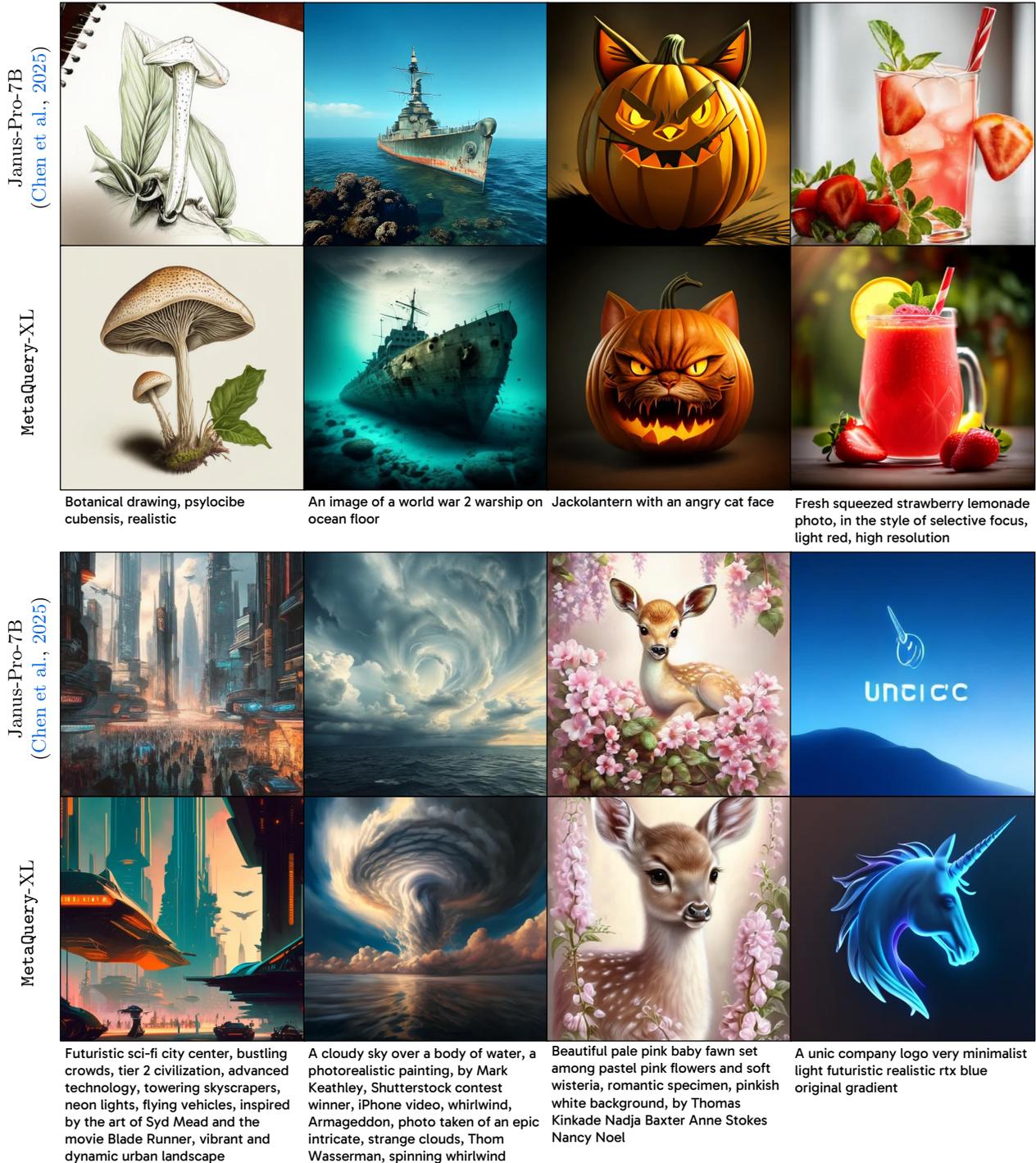
## C Training Objectives

Objective	Rel. Wall Time	MJHQ-30K FID↓	GenEval↑	DPG-Bench↑
Text-to-Image	1.0x	7.43	0.56	75.35
Image Reconstruction	2.79x	27.42	0.32	68.36
Mix	2.61x	8.27	0.54	76.53

**Table 10** Study on training objectives. Image reconstruction objective can be mixed with text-to-image objective to enable image reconstruction capabilities without harming visual quality and prompt alignment.

We are using an MLLM for multimodal perception, besides the standard text-to-image objective, we can also use an image reconstruction objective to achieve alignment. In Table 10, we show that training with the text-to-image objective achieves much better performance than the image reconstruction objective. We demonstrate that a mix of both objectives can enable image reconstruction capabilities without being generally harmful to the T2I performance.

<sup>2</sup><https://huggingface.co/Qwen/Qwen2-VL-7B-Instruct>



**Figure 10** Qualitative comparison with Janus-Pro-7B (Chen et al., 2025) on MJHQ-30K (Li et al., 2024b).