

SEED-X: Multimodal Models with Unified Multi-granularity Comprehension and Generation

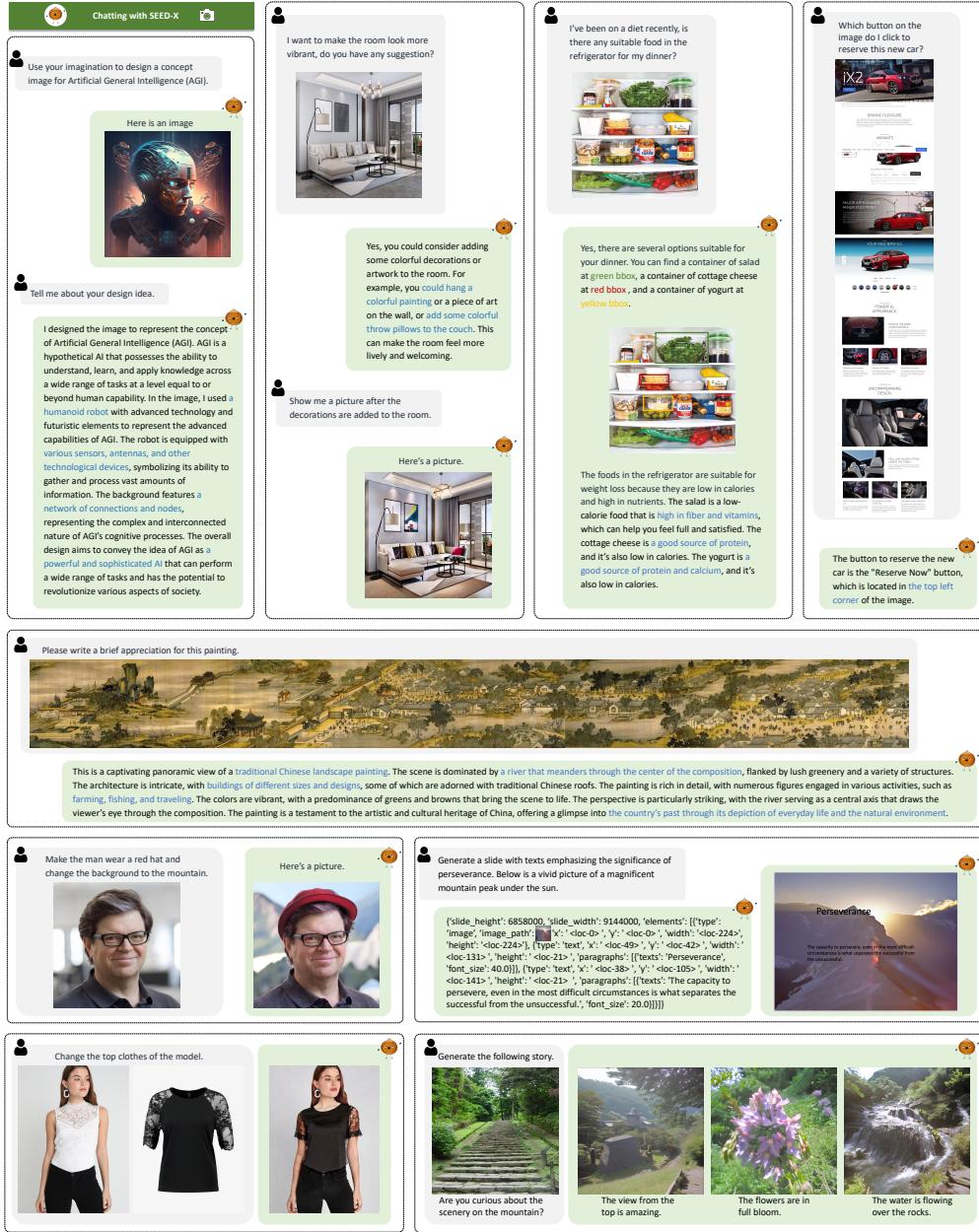
Yuying Ge^{1*}Sijie Zhao^{1*}Jinguo Zhu^{1*}Yixiao Ge^{1,2†}Kun Yi²Lin Song¹Chen Li²Xiaohan Ding¹Ying Shan^{1,2}¹Tencent AI Lab²ARC Lab, Tencent PCG

Figure 1: The introduced SEED-X, a unified and versatile foundation model, can serve as various multimodal AI assistants in the real world after different instruction tuning, capable of responding to a variety of user needs through unifying multi-granularity comprehension and generation.

Abstract

The rapid evolution of multimodal foundation model has demonstrated significant progresses in vision-language understanding and generation, *e.g.*, our previous work SEED-LLaMA. However, there remains a gap between its capability and the real-world applicability, primarily due to the model’s limited capacity to effectively respond to various user instructions and interact with diverse visual data. In this work, we focus on bridging this gap through integrating two enhanced features: (1) **comprehending images of arbitrary sizes and ratios**, and (2) **enabling multi-granularity image generation**. We present a unified and versatile foundation model, namely, **SEED-X**, which is able to model multi-granularity visual semantics for comprehension and generation tasks. Besides the competitive results on public benchmarks, SEED-X demonstrates its effectiveness in handling real-world applications across various domains after instruction tuning. We hope that our work will inspire future research into what can be achieved by versatile multimodal foundation models in real-world applications. The models, codes, and datasets are released in [https://github.com/AILab-CVC/SEED-X¹](https://github.com/AILab-CVC/SEED-X).

1 Introduction

In recent years, Multimodal Large Language Models (MLLMs) [1, 2, 3, 4, 5, 6, 7, 8] have demonstrated exceptional capabilities in comprehending multimodal data through leveraging the strong generality of LLMs [9, 10, 11]. Some pioneering work [12, 13, 14, 15, 16, 17, 18, 19] further empower LLMs with the ability to generate images beyond texts. For example, our previous work SEED-LLaMA [15] can handle a variety of tasks and excel in academic benchmarks through unifying multimodal comprehension and generation. However, the accuracy and diversity of its generated content still fall short of real-world needs. In this work, we focus on bridging this gap through upgrading SEED-LLaMA with enhanced capabilities for real-world applications.

Specifically, in order to make a multimodal foundation model applicable in real-world scenarios, we incorporate two enhanced features: (1) **understanding images of arbitrary sizes and ratios**, and (2) **multi-granularity image generation**, encompassing both **high-level instructional image generation** and **low-level image manipulation tasks**. These attributes can form the basis for a multimodal foundation model’s effective application in an open-world context, since a multimodal foundation model has to accommodate various downstream tasks requiring different levels of visual semantics.

In this paper, we introduce SEED-X, a unified and versatile multimodal foundation model as a follow-up work of SEED-LLaMA, which seamlessly integrates the features mentioned above. It is important to emphasize that *integrating all these characteristics into a single foundation model is by no means trivial*, as shown in Table 1, since none of the previous works support all of these features.

After different instruction tuning, SEED-X can function as various multimodal AI assistants in the real world, capable of addressing various user needs through generating proper texts and images as shown in Fig. 1. Specifically, our instruction-tuned models can act as an interactive designer, generating images while illustrating creative intent, offering modification suggestions and showcasing visualizations based on user’s input images. Additionally, they can act as knowledgeable personal assistants, comprehending images of various sizes and providing relevant suggestions. Moreover, they can generate more diverse outputs, such as slide layouts for slide creation, and interleaved image-text content for storytelling. SEED-X signifies a notable advancement towards a versatile agent for users in the real world.

To endow SEED-X with the aforementioned characteristics, our approach incorporates (1) a **visual tokenizer to unify image comprehension and generation**, where its **multi-granularity de-tokenization phase facilitates image generation and high-precision image manipulation**, and (2) an **MLLM with dynamic resolution image encoding to enable the comprehension of images with arbitrary sizes and**

¹This is the v2 version. We added benchmark results (without updating models) and ablation study.

*Equal Contribution.

[†]Correspondence to yixiaoge@tencent.com.

[‡]We sincerely acknowledge Tianheng Cheng (ARC Lab, Tencent PCG) for his support.

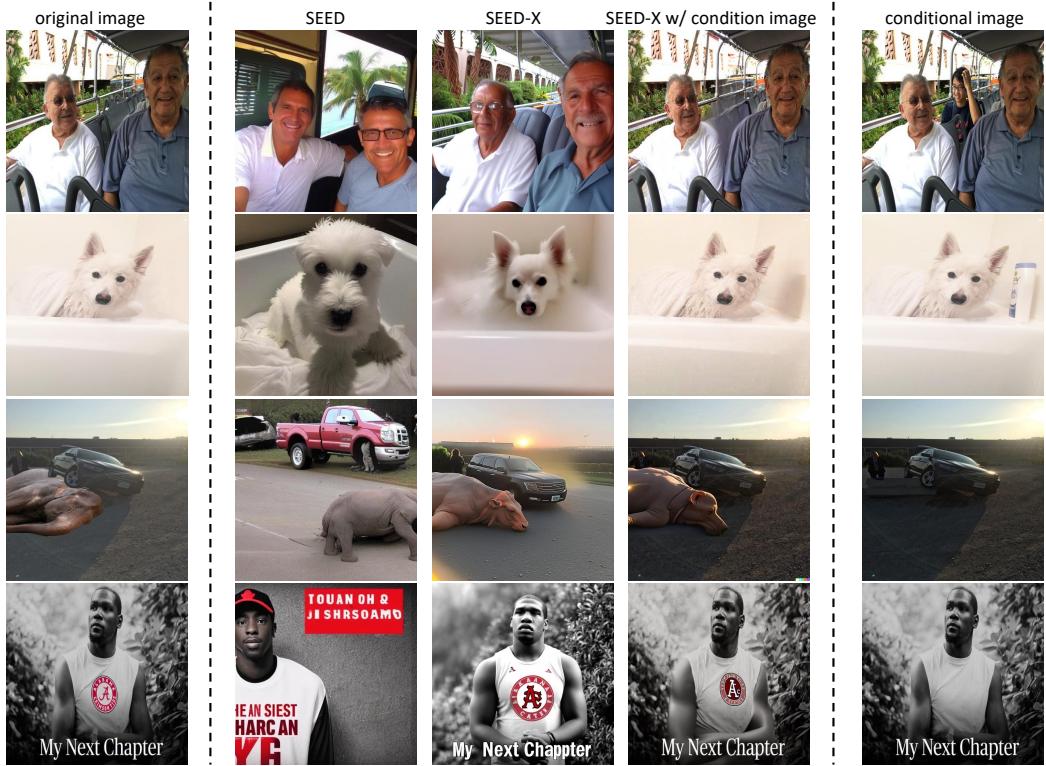


Figure 2: The reconstruction results of our visual de-tokenizer. It can decode realistic images that are semantically aligned with the original images by taking the **ViT features** as inputs, and further recover fine-grained details by incorporating the **conditional images** as inputs.

aspect ratios. Specifically, we utilize a **pre-trained ViT as the visual tokenizer** and train a **visual de-tokenizer** to decode realistic images by taking the **ViT features as input**. To realize the retention of fine-grained details of the input image to satisfy image manipulation, we further **fine-tune the visual de-tokenizer to take an extra condition image as input in the latent space** (See Fig. 2). The **ViT features serve as a bridge** to decouple the training of the visual (de-)tokenizer and the **MLLM**. The dynamic resolution image encoding divides an input image into sub-images and adds extrapolatable **2D positional embeddings** to the **ViT features of each sub-image**, allowing the **MLLM** to scale to any image resolution. For image generation, a fixed number of learnable queries are fed into the **MLLM**, where the **output hidden states** are trained to reconstruct the **ViT features of the target images**. During inference, the image de-tokenizer can take both the output features from the **MLLM** and the condition image provided by users as input, ensuring that the decoded image can possess high-level semantics that meet the multimodal instructions and retain the low-level details.

We pre-train SEED-X on massive multimodal data, including image-caption pairs, grounded image-text data, interleaved image-text data, OCR data, and pure texts. We further apply multimodal instruction tuning to align SEED-X with human instructions across various domains, utilizing both existing datasets and newly collected datasets that cover image editing, text-rich, grounded and referencing QA, and slide generation tasks. The extensive evaluations on MLLM benchmarks demonstrate that our instruction-tuned model not only achieves competitive performance in multimodal comprehension, but also exhibits excellent instruction-following capabilities for image generation.

All the models, codes, and datasets are made publicly available. We hope our work can bring insights to the community about the potential of multimodal foundation models in real-world scenarios through unifying multi-granularity comprehension and generation.

2 Related Work

With the rapid development of Multimodal Large Language Models (MLLM), recent studies have been working on unified MLLMs that are capable of **multimodal comprehension and generation** as shown in Tab. 1. Some work [15, 14, 13, 20, 21, 22, 23, 24] utilize a discrete visual tokenizer to

Table 1: MLLMs that unify comprehension and generation and whether they support significant characteristics essential for real-world applications. “Decoder Input” denotes the inputs for image generation, where “Features” means continuous features, “Token” represents discrete tokens, “Text” implies text prompts.

	Date	Decoder Input	Detection	Dynamic -Res Img Input	Image Gen	High-precision Editing	Open-source
Emu	07/2023	Feature	×	×	✓	×	✓
CM3Leon	07/2023	Token	×	×	✓	×	✗
SEED-OPT	07/2023	Token	×	×	✓	✗	✗
LaVIT	09/2023	Token	×	×	✓	✗	✓
NExT-GPT	09/2023	Feature	×	×	✓	✗	✓
DreamLLM	09/2023	Feature	×	×	✓	✗	✗
SEED-LLaMA	10/2023	Token	×	×	✓	✗	✓
VL-GPT	12/2023	Feature	×	×	✓	✗	✗
Gemini	12/2023	Token	×	-	✓	✗	✗
Emu2	12/2023	Feature	×	×	✓	✗	✓
Unified-IO 2	12/2023	Token	✓	×	✓	✗	✓
Mini-Gemini	03/2024	Text	×	×	✓	✗	✓
SEED-X	04/2024	Feature	✓	✓	✓	✓	✓

perform multimodal autoregression with a unified next-word-prediction objective or masked visual token prediction. Some research efforts [12, 25, 19] have delved into multimodal autoregression with continuous representations, where each image in the multimodal sequence is tokenized into embeddings via a visual encoder, and then interleaved with text tokens for autoregressive modeling. During inference, the regressed visual embeddings will be decoded into an image by a visual decoder. Additionally, some studies [17, 16] enable image generation in a non-autoregressive manner through utilizing learnable queries to obtain visual representations from MLLMs, which are further fed into a image decoder to generate images. Mini-Gemini, generates text prompts using MLLMs and then leverages the existing SDXL [26] to output images.

Although these work have achieved competitive results on various academic benchmarks, such as VQA and text-to-image generation, the accuracy and diversity of their generated content still fall short of real-world needs, since they do not meet the requirements of modeling multi-granularity visual semantics for comprehension and generation task. As shown in Tab. 1, we identify several significant characteristics essential for real-world applications including object detection and dynamic resolution image encoding for multi-granularity comprehension, as well as high-level instructional image generation and low-level image manipulation for multi-granularity image generation. Notably, **none of the previous works fully support all of these characteristics**. In this work, we present SEED-X, a unified and versatile foundation model, which effectively incorporate the aforementioned characteristics for real-world applications.

3 Method

3.1 Visual Tokenization and De-tokenization

In SEED-X, we adopt a visual tokenizer to unify image comprehension and generation, and pre-train a multi-granularity de-tokenizer to facilitate image generation and high-precision image manipulation in a two-stage manner. In the first stage, as shown in Fig. 3 (left), we utilize a pre-trained ViT as the visual tokenizer and pre-train a visual de-tokenizer to decode realistic images by taking the features of the ViT as inputs in the first stage. Specifically, N visual embeddings from the ViT tokenizer ($N = 64$ after average pooling) are fed into a learnable module as the inputs of the U-Net of the pre-trained SD-XL [26] (replacing the original text features). The learnable module consists of four cross-attention layers to connect the visual tokenizer and the U-Net. We optimize the parameters of the learnable module and keys and values within the U-Net on the images from JourneyDB [27], LAION-Aesthetics [28], Unsplash [29], and LAION-COCO [30]. As shown in Fig. 2, compared with

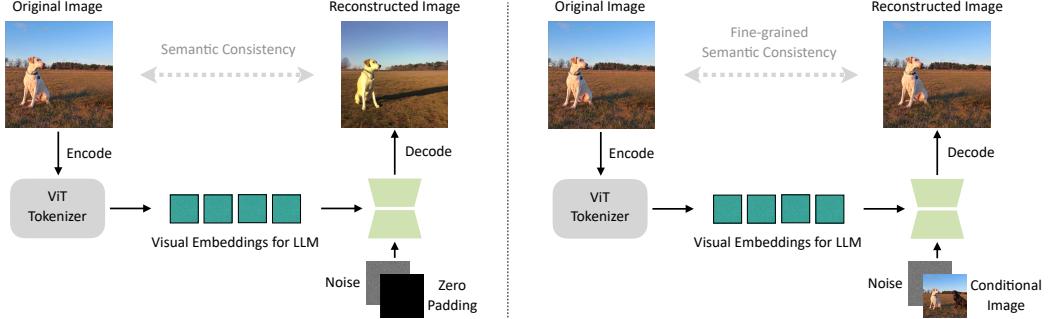


Figure 3: Overview of visual tokenization and de-tokenization in SEED-X. In the first stage (left), we pre-train a visual de-tokenizer, which can decode semantically consistent images by taking the features of a pre-trained ViT as inputs. In the second stage (right), we fine-tune the visual de-tokenizer through concatenating the latent features of a conditional image with the noise to recover the fine-grained details of the original image.

SEED [15], our visual de-tokenizer can decode images that are more semantically aligned with the original images by taking the ViT features as inputs.

In the second stage, as shown in Fig. 3 (right), we further fine-tune the visual de-tokenizer to take an extra condition image as inputs for the retention of low-level details. Specifically, we follow InstructPix2Pix [31] to encode the condition image into the latent space via the VAE encoder, and concatenate them with the noisy latent as the input of U-Net. The channel number of the U-Net convolutional layer is expanded from 4 to 8, and all parameters of U-Net are optimized. We fine-tune the visual de-tokenizer on MagicBrush [32] and in-house image editing data, as well as the pure images in the first stage, where the conditional inputs are set to zeros. As shown in Fig. 2, by incorporating the condition image as an additional input besides the high-level image features, our visual de-tokenizer can recover the fine-grained details of the original image.

3.2 Dynamic Resolution Image Encoding

Current MLLMs require to resize the input images to a pre-defined resolution (typically a square size), which corresponds to the training resolution of the vision encoder, which can result in the loss of fine-grained information. In this work, we propose dynamic resolution image encoding to enable the processing of images with arbitrary sizes and aspect ratios by dividing the image into a grid comprising of sub-images. Specifically, for the visual encoder with the training resolution $H_t \times W_t$, we first up-sample the input image with the size $H \times W$ to the size of $\{N_h * H_t\} \times \{N_w * W_t\}$. The grid size $N_h \times N_w$, are determined by

$$\begin{aligned} \min \quad & N_h * N_w, \\ \text{s.t. } & H \leq N_h * H_t \quad \text{and} \quad W \leq N_w * W_t. \end{aligned} \tag{1}$$

We also resize the original image to the size of $H_t \times W_t$ to provide global visual context. All sub-images and the resized global image are fed into the visual encoder to obtain the features, which are concatenated as the input of the LLM.

To enable the LLM to be aware of the positional information of each sub-image within the original image, we add extrapolatable 2D positional embeddings to the visual features of each sub-image. Specifically, for a sub-image with a normalized center location (x_c, y_c) in the grid, where $0.0 < x_c, y_c < 1.0$, its learnable positional embedding p is computed:

$$p = x_c * l + (1 - x_c) * r + y_c * t + (1 - y_c) * b. \tag{2}$$

l , r , t , and b represent four learnable position embeddings indicating left, right, top and bottom respectively. Consequently, our visual encoder can handle inputs with any arbitrary sizes and aspect ratios, even if the image resolution was not encountered during training.

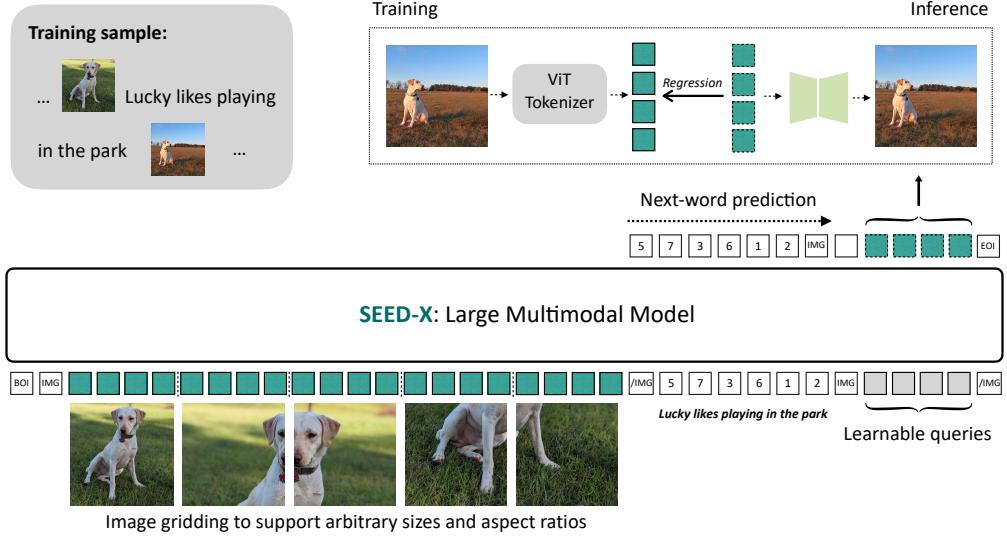


Figure 4: Overview of SEED-X for multimodal pre-training. Each image is divided into sub-images to support arbitrary sizes and aspect ratios, and their ViT features along with text tokens are fed into an LLM to perform next-word prediction and image feature regression between the output hidden states of the learnable queries and ViT features. During inference, the regressed image features are fed into the visual de-tokenizer to decode images.

3.3 Multimodal Pre-training and Instruction Tuning

3.3.1 Training Stage I: Multimodal Pre-training

As shown in Fig. 4, SEED-X adopts next-word prediction and image feature regression training objectives on interleaved visual and textual data. Specifically, we perform dynamic resolution encoding of each image in the multimodal sequence, and their features along with text tokens are fed into the pretrained LLM. In order to equip the model with detection and referencing abilities, we add **224 bbox tokens, designated for representing bounding box coordinates, represented by <box_start> <loc-x_center> <loc-y_center> <loc-width> <loc-height> <box_end>** with special tokens at the beginning and end of the bounding box. The text and added bbox tokens are trained through predicting the next token with cross-entropy loss.

We employ N learnable queries ($N = 64$ to align with the visual de-tokenizer) to obtain the output visual representations from the LLM, which are trained to **reconstruct the features of the pre-trained ViT tokenizer with a Mean Squared Error (MSE) loss**. We add two special tokens ‘’ and ‘’ to represent the **beginning and the end of the query embeddings**, and the ‘’ is trained to predict where an image emerges. In doing so, we utilize the pre-trained ViT tokenizer as a **bridge** to decouple the training of a visual de-tokenizer and the MLLM for image generation. During inference, the regressed visual representations from SEED-X are fed into the visual de-tokenizer to decode realistic images.

We pre-train SEED-X initialized from **Llama2-chat-13B** using LoRA on massive multimodal data, including image-captions pairs, grounded image-texts, interleaved image-text data, OCR data and pure texts. We perform pre-training with **48 H800-80G GPUs (10 days)** on a total of 158M samples. See Appendix. A and Appendix. B for more details.

3.3.2 Training Stage II: Multimodal Instruction Tuning

We perform multimodal instruction tuning through fine-tuning SEED-X using a LoRA module with both public datasets and in-house data covering image editing, text-rich, grounded and referencing QA, and slide generation tasks. The details of datasets can be found in Appendix. A. We fine-tune SEED-X with conversational and image generation data to yield a general instruction-tuned model SEED-X-I, which can follow multimodal instructions and make responses with images, texts

Table 2: **Comparison on multimodal understanding benchmarks.** “Und.” and “Gen.” denote “understanding” and “generation”, respectively.

Type	Model	POPE↑	MME-P↑	MMB↑	SEED _(img) ↑	VQAv2 _(test) ↑	GQA↑	MMMU↑	MM-Vet↑
Und. Only	LLaVA-v1.5-Phi-1.5 [23]	84.1	1128.0	-	-	75.3	56.5	30.7	-
	MobileVLM [33]	84.5	1196.2	53.2	-	-	56.1	-	-
	MobileVLM-V2 [34]	84.3	1302.8	57.7	-	-	59.3	-	-
	LLaVA-Phi [35]	85.0	1335.1	59.8	-	71.4	-	-	28.9
	LLaVA [36]	76.3	809.6	38.7	33.5	-	-	-	25.5
	LLaVA-v1.5 [37]	85.9	1510.7	64.3	58.6	78.5	62.0	35.4	31.1
	InstructBLIP [38]	-	-	36.0	53.4	-	49.2	-	26.2
	IDEFICS-9B [39]	-	-	48.2	-	50.9	38.4	-	-
Und. and Gen.	Qwen-VL-Chat [5]	-	1487.5	60.6	58.2	78.2	57.5	-	-
	DreamLLM [17]	-	-	-	-	72.9	-	-	36.6
	LaVIT [20]	-	-	-	-	66.0	46.8	-	-
	Emu [18]	-	-	-	-	52.0	-	-	-
	NExT-GPT [40]	-	-	-	-	66.7	-	-	-
	Gemini-Nano-1 [41]	-	-	-	-	62.7	-	26.3	-
	LWM [42]	75.2	-	-	-	55.8	44.8	-	9.6
	SEED-X	84.1	1457.0	70.1	66.5	71.2	49.1	35.6	43.0

Table 3: **Evaluation of text-to-image generation ability on GenEval benchmark.** “Und.” and “Gen.” denote “understanding” and “generation”, respectively.

Type	Method	Single Obj.	Two Obj.	Counting	Colors	Position	Color Attri.	Overall↑
Gen. Only	LDM [52]	0.92	0.29	0.23	0.70	0.02	0.05	0.37
	SDv1.5 [52]	0.97	0.38	0.35	0.76	0.04	0.06	0.43
	PixArt- α [53]	0.98	0.50	0.44	0.80	0.08	0.07	0.48
	SDv2.1 [52]	0.98	0.51	0.44	0.85	0.07	0.17	0.50
	DALL-E 2 [54]	0.94	0.66	0.49	0.77	0.10	0.19	0.52
	SDXL [26]	0.98	0.74	0.39	0.85	0.15	0.23	0.55
Und. and Gen.	LWM [42]	0.93	0.41	0.46	0.79	0.09	0.15	0.47
	SEED-X	0.96	0.65	0.31	0.80	0.18	0.14	0.51

and bounding boxes in multi-turn conversation. We further fine-tune the foundation model SEED-X on specialized datasets, resulting in a series of instruction-tuned models tailored for specific tasks, including SEED-X-Edit, SEED-X-PPT, SEED-X-Story and SEED-X-Try-on. The proficient capabilities of these instruction-tuned model across various domains demonstrate the versatility of our pre-trained foundation model SEED-X. We perform instruction tuning on the foundation model SEED-X across different datasets, resulting in various models with distinct capabilities. Note that **we do not have an all-in-one instruction-tuned model that encompasses all abilities**, which will be explored for future work.

4 Experiments

4.1 Quantitative Evaluation

Multimodal Comprehension. We evaluate the multimodal comprehension capabilities of SEED-X-I on widely recognized image-based vision-language benchmarks, which include VQAv2 [43], GQA [44], POPE [45], MME [46], SEED [47], MMB [48], MM-Vet [49], and MMMU [50]. As listed in Tab. 2, SEED-X-I achieves competitive performance across various benchmarks, even when compared to MLLMs specifically designed for multimodal comprehension.

Image Generation. We evaluate the image generation capabilities of SEED-X-I on GenEval [51], which is a challenging benchmark to evaluate compositional image properties such as object co-occurrence, position, count, and color. As shown in Tab. 3, SEED-X obtains 51% overall accuracy, demonstrating the model’s excellent instruction-following capabilities for image generation.

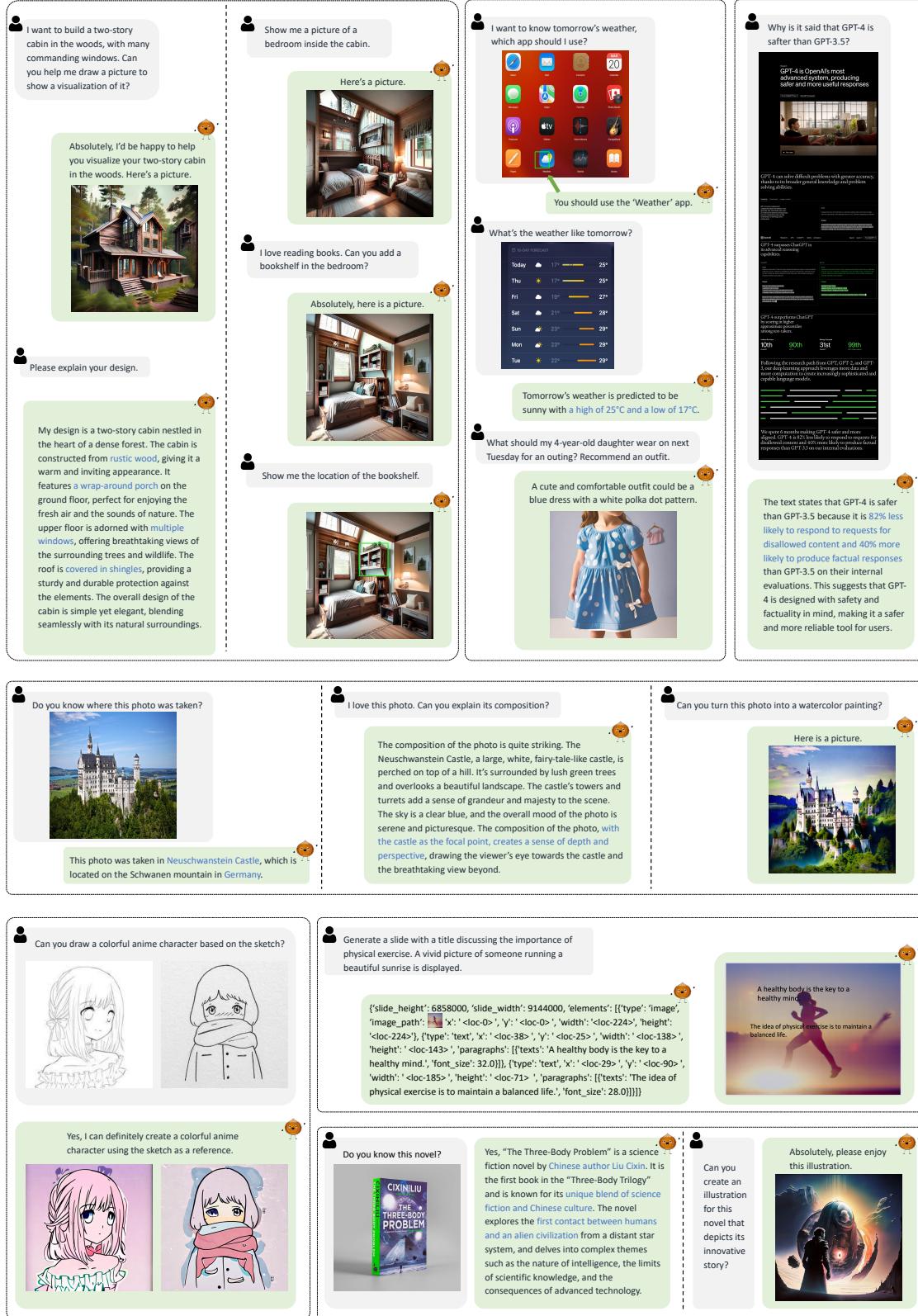


Figure 5: Examples of what SEED-X can do in real-world scenarios after different instruction tuning through unifying multi-granularity comprehension and generation. Our instruction tuned models can function as an interactive designer, generating images without descriptive captions while illustrating creative intent, and showcasing visualizations of modified images based on user’s intent. They can act as knowledgeable personal assistants, comprehending images of arbitrary sizes and offering relevant suggestions in multi-turn conversations.

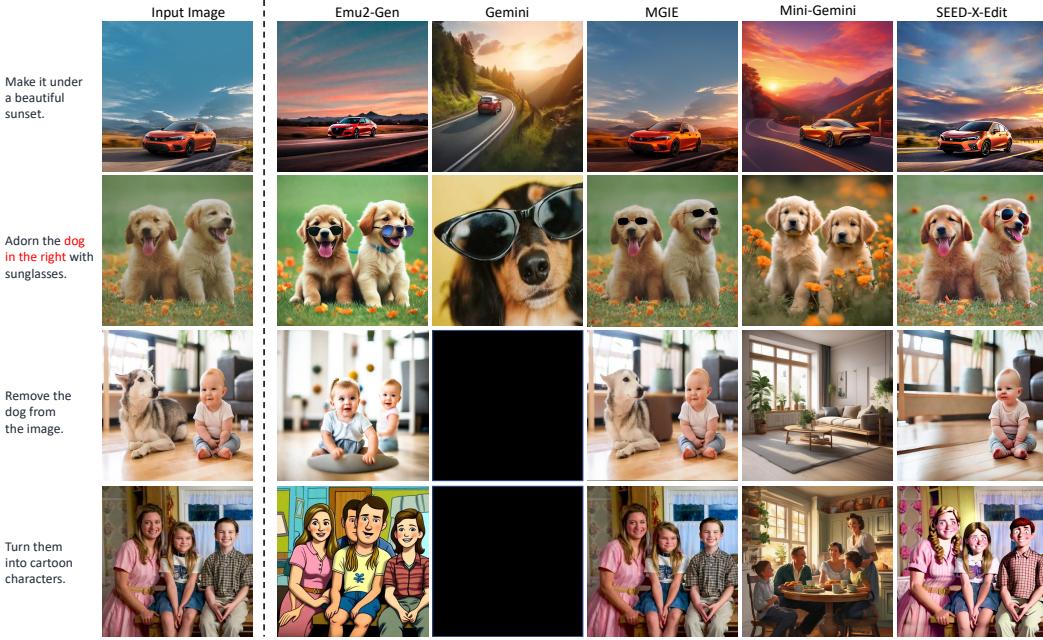


Figure 6: Qualitative comparison between MLLMs for image manipulation. SEED-X-Edit shows enhanced ability in adhering to instructions while preserving low-level details of input images. The black images result from Gemini’s inability to display human images.

4.2 Qualitative Evaluation

4.2.1 Applications in the Real World.

SEED-X can be effectively instruction tuned to function as various multimodal AI assistants in the real world across different domains after integrating two enhanced features, including the comprehension of images of arbitrary sizes and ratios, and multi-granularity image generation, encompassing both high-level instructional image generation and low-level image manipulation tasks. As shown in Fig. 1 and Fig. 5, our instruction tuned models can serve as an interactive designer, which can generate images without descriptive captions while illustrate creative intent, and showcase visualizations of modified images. For example, it can explain the design idea of concept image for AGI and a two-story cabin. It can create an imaginative illustration for the novel without the need of describing the scene with languages. It can further offer modification suggestions of the user’s room and showcase the visualization. Additionally, the instruction tuned models can act as an knowledgeable personal assistant, comprehending images of arbitrary sizes and providing relevant suggestions. For example, it can identify foods suitable for fat reduction in the refrigerator, display appropriate clothing based on the screenshot of weather forecasts.

4.2.2 Image Generation and Manipulation.

We compare previous MLLMs that are capable of generating images for text-to-image generation in Fig. 9 of Appendix. Our instruction tuned model can generate images that are more aligned with the elements in the caption and possess artistic qualities. Through utilizing a pre-trained ViT Tokenizer as the bridge to decouple the training of visual de-tokenizer and the MLLM, our pre-trained model SEED-X can effectively realize high-quality image generation, which is a fundamental capability to be applied in real-world scenarios.

We compare image manipulation with previous MLLMs including Emu2-Gen [25], Gemini [41], MGIE [41] and Mini-Gemini [55]. As shown in Fig. 6, we can observe that SEED-X-Edit can more effectively adhere to editing instructions while maintaining the low-level details of the input image. For instance, SEED-X-Edit can accurately add sunglasses to the dog on the right, while both Emu2-Gen and MGIE fail to follow the instruction, resulting in sunglasses being added to both dogs.

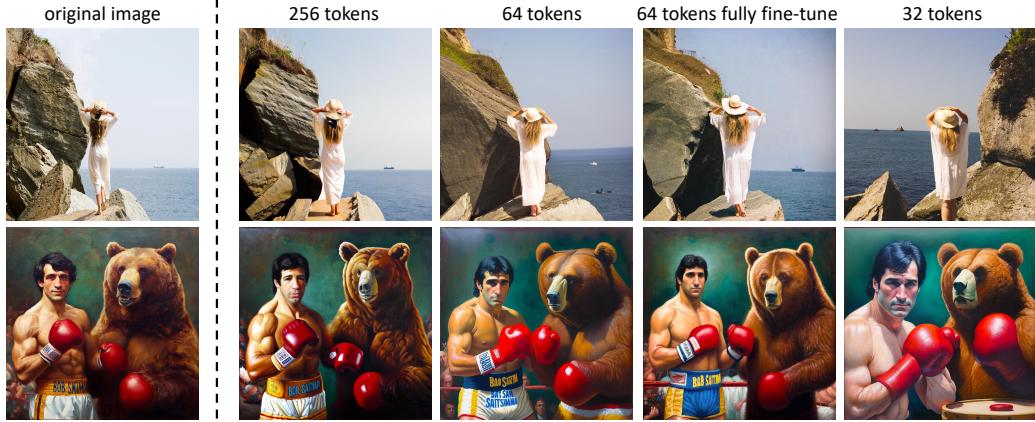


Figure 7: Ablation study on the number of visual tokens and trainable parameters for training visual de-tokenizer.

Additionally, SEED-X-Edit successfully eliminates the dog in the baby image while preserving the background details and the baby’s features. In contrast, Emu2-Gen fails to retain the fine details of the input image, and MGIE is unsuccessful in removing the dog. Note that Gemini lacks the ability to edit images as it **retrieves images** on the Internet. Here the presence of black images is due to its failure to display images related to human portraits. Mini-Gemini generates **text prompts** as the input of a pre-trained SDXL model, which can not preserve the visual details of the input image. The examples show the effectiveness of our instruction model for high-precision image manipulation. Our MLLM accurately predicts visual semantic representations based on an input image and a language instruction, which serve as input for the U-Net. The visual de-tokenizer can further condition on the input image, ensuring the preservation of fine-grained details in the decoded images.

4.2.3 Multimodal Comprehension.

We provide qualitative examples of multimodal comprehension by SEED-X-I in Fig. 10 and Fig. 11 of Appendix. SEED-X-I can realize fine-grained object detection and perception, text-rich comprehension, fundamental mathematical computation, world-knowledge and commonsense reasoning, diagram understanding, which are crucial capabilities for its application in real-world scenarios.

4.3 Ablation Study

In this section, we perform ablation studies on the training of our visual de-tokenizer and the pre-training of SEED-X to enable a MLLM for image generation.

For visual de-tokenization, N visual embeddings (after average pooling) from the ViT tokenizer are fed into a learnable module as the inputs of the U-Net of the pre-trained SD-XL. We perform an ablation study on the number of visual tokens and the learnable parameters of the SD-XL U-Net, where keys and values within the U-Net are optimized if not specified with “fully fine-tune”. As shown in Fig. 7, we can observe that more visual tokens can result in better reconstruction of the original images. For example, the decoded images from 256 visual embeddings can recover the characters’ postures of the original images, while decoded images from 32 visual embeddings have already lost the original structure of the scene. We further observe that **fully fine-tuning the parameters of the SD-XL U-Net can lead to distortions in image details**, such as the woman’s feet, **compared to only training the keys and values within the U-Net**. In SEED-X, we use $N = 64$ visual embeddings to train the visual de-tokenizer and only optimize the keys and values within the U-Net (See below for an explanation of why we do not choose $N = 256$).

To enable MLLM for image generation, we employ N learnable queries to obtain the output visual representations from the LLM, which are trained to reconstruct N visual embeddings from the ViT tokenizer with a learnable module. We first perform an ablation study on the number of learnable queries. The images generated by the MLLM based on the input caption are shown in Fig. 8. We can observe that using 256 learnable queries to reconstruct 256 visual embeddings can lead to distortion

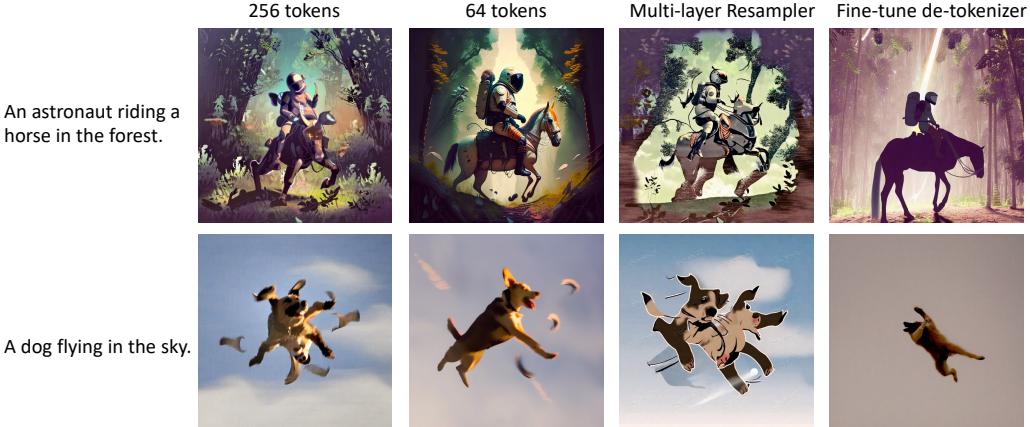


Figure 8: Ablation study on the number of visual tokens, model architecture and optimization targets during pre-training SEED-X for image generation.

in the generated images compared with $N = 64$. This occurs because regressing more visual features is more challenging for the model, even though 256 visual embeddings from the de-tokenizer can better reconstruct images, as demonstrated in the previous ablation study. We also observe that, compared to learning a one-layer cross-attention for reconstructing image features, a multi-layer resampler (multi-layer cross-attention) yields less satisfactory performance, which can happen due to the lack of more direct regularizations on the hidden states of the LLM. We further optimize the visual de-tokenizer by using the reconstructed visual embeddings from the MLLM as input instead of ViT features, but the generated images exhibit a more monotonous appearance. It demonstrates the effectiveness of utilizing the ViT Tokenizer as the bridge to decouple the training of visual de-tokenizer and the MLLM for image generation.

5 Conclusion

We present SEED-X, a versatile foundation model, which can serve as various multimodal AI assistants in the real world after instruction tuning. In order to make a multimodal foundation model applicable in open-world context, we integrate two enhanced features into SEED-X including image comprehension of arbitrary sizes and ratios, and multi-granularity image generation, which encompasses both high-level instructional image generation and low-level image manipulation. We hope that SEED-X can inspire future research into the potential of MLLMs in the real-world scenarios through unifying multi-granularity comprehension and generation.

References

- [1] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICML*, 2023.
- [2] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [3] Haotian Liu, Chunyuan Li, Qingsyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- [4] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [5] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [6] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.

- [7] Pan Zhang, Xiaoyi Dong, Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*, 2023.
- [8] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023.
- [9] 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.
- [10] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [12] Quan Sun, Qiying 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.
- [13] Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: Pretraining and instruction tuning. *arXiv preprint arXiv:2309.02591*, 2023.
- [14] 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.
- [15] Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer. *arXiv preprint arXiv:2310.01218*, 2023.
- [16] 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.
- [17] Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. *arXiv preprint arXiv:2309.11499*, 2023.
- [18] Quan Sun, Qiying 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.
- [19] Jinguo Zhu, Xiaohan Ding, Yixiao Ge, Yuying Ge, Sijie Zhao, Hengshuang Zhao, Xiaohua Wang, and Ying Shan. Vl-gpt: A generative pre-trained transformer for vision and language understanding and generation. *arXiv preprint arXiv:2312.09251*, 2023.
- [20] 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. *arXiv preprint arXiv:2309.04669*, 2023.
- [21] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action. *arXiv preprint arXiv:2312.17172*, 2023.
- [22] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- [23] 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.
- [24] 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.
- [25] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, et al. Generative multimodal models are in-context learners. *arXiv preprint arXiv:2312.13286*, 2023.
- [26] 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.

- [27] Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances in Neural Information Processing Systems*, 36, 2024.
- [28] Christoph Schuhmann and Romain Beaumont. Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics/>, 2022.
- [29] Zahid Ali, Chesser Luke, and Carbone Timothy. Unsplash. <https://github.com/unsplash/datasets>, 2023.
- [30] Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. Laion-coco: 600m synthetic captions from laion2b-en. <https://laion.ai/blog/laion-coco/>, 2023.
- [31] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- [32] Kai Zhang, Lingbo Mo, Wenhua Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *arXiv preprint arXiv:2306.10012*, 2023.
- [33] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023.
- [34] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*, 2024.
- [35] Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. Llava-phi: Efficient multi-modal assistant with small language model. *arXiv preprint arXiv:2401.02330*, 2024.
- [36] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [37] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [38] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- [39] Hugo Laurençon, Daniel van Strien, Stas Bekman, Leo Tronchon, Lucile Saulnier, Thomas Wang, Siddharth Karamcheti, Amanpreet Singh, Giada Pistilli, Yacine Jernite, and et al. Introducing idefics: An open reproduction of state-of-the-art visual language model, 2023.
- [40] 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.
- [41] Gemini Team, Rohan Anil, Sébastien Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [42] Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024.
- [43] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [44] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [45] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- [46] 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.
- [47] 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*, 2023.
- [48] 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.

- [49] 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.
- [50] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.
- [51] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [53] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- [54] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [55] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [56] Schuhmann Christoph, Köpf Andreas, Vencu Richard, Coombes Theo, and Beaumont Romain. Laion coco: 600m synthetic captions from laion2b-en. [EB/OL], 2022. <https://laion.ai/blog/laion-coco/>.
- [57] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [58] Junting Pan, Keqiang Sun, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *arXiv preprint arXiv:2307.00716*, 2023.
- [59] Qiying Yu, Quan Sun, Xiaosong Zhang, Yufeng Cui, Fan Zhang, Xinlong Wang, and Jingjing Liu. Capsfusion: Rethinking image-text data at scale. *arXiv preprint arXiv:2310.20550*, 2023.
- [60] 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. *arXiv preprint arXiv:2304.06939*, 2023.
- [61] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text documents, 2023.
- [62] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.
- [63] Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.
- [64] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*, 2023.
- [65] Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*, 2019.
- [66] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- [67] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 235–251. Springer, 2016.
- [68] Pan Lu, Swaroop Mishra, Tanglin 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. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.
- [69] Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. Kvqa: Knowledge-aware visual question answering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8876–8884, 2019.

- [70] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2018.
- [71] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- [72] Chen Li, Yixiao Ge, Dian Li, and Ying Shan. Vision-language instruction tuning: A review and analysis. *arXiv preprint arXiv:2311.08172*, 2023.
- [73] Junke Wang, Lingchen Meng, Zejia Weng, Bo He, Zuxuan Wu, and Yu-Gang Jiang. To see is to believe: Prompting gpt-4v for better visual instruction tuning. *arXiv preprint arXiv:2311.07574*, 2023.
- [74] Zhiyang Xu, Chao Feng, Rulin Shao, Trevor Ashby, Ying Shen, Di Jin, Yu Cheng, Qifan Wang, and Lifu Huang. Vision-flan: Scaling human-labeled tasks in visual instruction tuning. *arXiv preprint arXiv:2402.11690*, 2024.
- [75] 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, 2024.
- [76] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallochi, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International journal of computer vision*, 128(7):1956–1981, 2020.
- [77] Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239, 2016.
- [78] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. Viton-hd: High-resolution virtual try-on via misalignment-aware normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14131–14140, 2021.

Table 4: Overview of the pre-training and instruction tuning datasets.

Type	Dataset
Pre-training	
Image-Caption	LAION-COCO [56] (Re-caption), SAM [57] (Re-caption), Unsplash [29], LAION-Aesthetics[28], JourneyDB [58], CapFusion [59],
Grounded Image-Caption	GRIT [4]
Interleaved Image-Text	MMC4 [60], OBELICS [61], OpenFlamingo [62]
OCR	LLaVAR [63], Slides (In-house)
Pure Text	Wikipedi
Instruction Tuning	
VQA	LLaVAR [63], Text-rich QA (In-house), MIMIC-IT [64], MathQA [65], ChartQA [66], AI2D [67], ScienceQA [68], KVQA [69], DVQA [70], Grounded QA (In-house), Referencing QA (In-house)
Conversation	LLaVA-150k [36], ShareGPT [71], VLIT [72], LVIS-Instruct4V [73], Vision-Flan [74], ALLaVA-4V [75]
Image Generation	LAION-COCO [56] (Re-caption), SAM [57] (Re-caption), Unsplash [29], LAION-Aesthetics[28], JourneyDB [58]
Image Editing	Instructpix2pix [31], MagicBrush [32], Openimages [76]-editing (In-house), Unsplash [29]-editing (In-house)
Slides Generation	In-house data
Story Telling	VIST [77]
Virtual Try-on	VITON-HD [78]

A Pre-training and Instruction Tuning Datasets

As listed in Tab. 4, we pre-train SEED-X and conduct instruction tuning on a large variety of both public datasets and in-house data. For multimodal pre-training, we utilize image-caption pairs, grounded image-caption pairs, interleaved image and text content, OCR data and pure text data. The images of LAION-COCO [56] and SAM [57] are re-captioned for a more detailed descriptive caption to improve both image comprehension and generation.

For instruction tuning, we utilize various public VQA datasets, and further curate text-rich QA, grounded and referencing QA to enhance the model’s capability of comprehending text-rich images and detecting objects that requires reasoning. We use multiple conversational datasets, which are specifically collected for MLLMs with open-form text output. We use the same image-caption pairs as in the pre-training phase to maintain the model’s ability to generate images. For the image manipulation, since the high-precision editing dataset MagicBrush [32] is only at the level of thousands, we employ a series of models to collect a dataset of millions of image editing examples, which are used for both training the visual de-tokenizer and SEED-X-Edit. We further collected data on slides, obtaining images, captions, and layouts for training slide generation.

B Implementation Details

Visual Tokenization and De-tokenization. We use the visual encoder from Qwen-vl [5] as the ViT Tokenizer and adopt 1D average pooling to obtain $N = 64$ visual embeddings. These visual embeddings are fed into four layers of cross-attention as the input of the U-Net initialized from SDXL [26]. In the first stage, we optimize the parameters of the cross-attention layers and the keys and values within the U-Net on the images from JourneyDB [27], LAION-Aesthetics [28], Unsplash [29], and LAION-COCO [30]. We train the visual de-tokenizer on 32 A100-40G GPUs with 42K training steps, where the learning rate is set to 1e-4 with cosine decay.

In the second stage, we encode the condition image into the latent space via the VAE encoder, and concatenate them with the noisy latent as the input of U-Net. The channel number of the U-Net convolutional layer is expanded from 4 to 8, and all parameters of U-Net are optimized. We pre-train the visual conditioner on MagicBrush [32] and in-house image editing data, as well as the image-caption pairs in the first stage, where the conditional inputs are set to zeros. We fine-tune the visual de-tokenizer on 32 A100-40G GPUs with 30K training steps, where the learning rate is set to 1e-4 with cosine decay.

Multimodal Pre-training and Instruction Tuning. We utilize the [visual encoder from Qwen-vl](#) [5] as the ViT Tokenizer and initialize a cross-attention layer to obtain $N = 64$ visual embedding as the input of the LLM initialized from Llama2-chat-13B. We initialize $N = 64$ learnable queries and the output hidden states from them are fed into a cross-attention layer to reconstruct $N = 64$ visual embeddings from the ViT Tokenizer. We optimize the LLM using LoRA and optimize the parameters of the input cross-attention layer, output cross-attention layer, extrapolatable 2D positional embeddings, and LoRA on image-captions pairs, grounded image-texts, interleaved image-text data, OCR data and pure texts. We perform pre-training with 48 H800-80G GPUs (10 days) on a total of 158M samples, where the learning rate is set to 1e-4 with cosine decay.

For the instruction tuning, we fine-tune a LoRA module on the pre-trained model, and optimize the parameters of the input cross-attention layer, output cross-attention layer, extrapolatable 2D positional embeddings, and LoRA. We fine-tune SEED-X with conversational and image generation data to yield a general instruction-tuned model SEED-X-I. We further fine-tune SEED-X on specialized datasets, resulting in a series of instruction-tuned models tailored for specific tasks, including SEED-X-Edit, SEED-X-PPT, SEED-X-Story and SEED-X-Try-on.

C Qualitative Examples

Text-to-image Generation. Fig. 9 visualizes the comparison between MLLMs for text-to-image generation including Next-GPT [16], SEED-LLaMA-I[15], Emu2-Gen [25] and Gemini [41]. Compared with previous MLLMs, our instruction tuned model can generate images that are more aligned with the elements in the descriptive caption and possess artistic qualities. For example, images generated by SEED-X-I vividly and accurately depicts “person standing in a small boat”, “a gleaming sword on its back”, “an oriental landscape painting”, “tiger with vivid colors” in the captions. Through utilizing a pre-trained ViT Tokenizer as the bridge to decouple the training of visual de-tokenizer and the MLLM, our pre-trained model SEED-X can effectively realize high-quality image generation, which is a fundamental capability for applying multimodal models in real-world scenarios.

Image Manipulation. We compare image manipulation with previous MLLMs including Emu2-Gen [25], Gemini [41], MGIE [41] and Mini-Gemini [55]. Language-guided image manipulation presents a significant challenge as the model must be capable of comprehending free-form instructions and generating images with the low-level details of the input image preserved. As shown in Fig. 6, we can observe that SEED-X-Edit can more effectively adhere to editing instructions while maintaining the low-level details of the input image. For instance, SEED-X-Edit can accurately add sunglasses to the dog on the right, while both Emu2-Gen and MGIE fail to follow the instruction, resulting in sunglasses being added to both dogs. Additionally, SEED-X-Edit successfully eliminates the dog in the baby image while preserving the low-level background details and the baby’s features. In contrast, Emu2-Gen fails to retain the fine details of the input image, and MGIE is unsuccessful in removing the dog. Note that Gemini lacks the ability to edit images as it **retrieves images** on the Internet. Here the presence of black images is due to its failure to display images related to human portraits. Mini-Gemini generates **text prompts** as the input of a pre-trained SDXL model, which can not preserve the visual details of the input image. The examples demonstrate the effectiveness of our instruction model for high-precision image manipulation. Our MLLM accurately predicts visual semantic representations based on an input image and a language instruction, which serve as input for the U-Net. The visual de-tokenizer can further condition on the input image, ensuring the preservation of fine-grained details in the decoded images.

Multimodal Comprehension We show qualitative examples of multimodal comprehension by SEED-X-I in Fig. 10 and Fig. 11. SEED-X-I can realize fine-grained object detection and perception, text-rich comprehension, fundamental mathematical computation, world-knowledge and commonsense reasoning, diagram understanding, which are crucial capabilities for its application in the real world.



A lonely person standing in a small boat, floating in the vast ocean, surrounded by thick fog in the sky, giving a sense of confusion and helplessness.



A fearless cat, with a gleaming sword on its back.



The golden moment of sunrise, Huangshan, China, towering Qishi Peak, a large area of clouds, a small welcome pine, an oriental landscape painting.



A fierce cartoon tiger, with furrowed brows and bared teeth, displays its anger through vivid colors and exaggerated features.

Figure 9: Qualitative comparison between MLLMs for text-to-image generation. SEED-X-I is capable of generating images that are more closely aligned with the elements in the descriptive caption and exhibit artistic qualities.



Detect the boat in the image.



What is the size of the second bedroom?



The second bedroom is 10'10" x 14'7".



Here are the order details for my taxi ride. Should I actually pay \$29.42?

Order Details	
• Base Fare	\$10
• Mileage Fee	\$12.11
• Duration Fee	\$10.32
• Total	\$32.42
• Discount	\$3
• Actually Paid	\$29.42



Based on the image, the total amount due is \$32.42, which is the sum of the base fare, mileage fee, duration fee. The discount is \$3, which is subtracted from the total, resulting in the actual paid amount of \$29.42. Therefore, you should pay \$29.42 for the taxi ride.



Count the number of the people.



There are 11 people in the image.



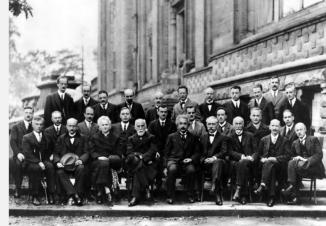
What are all the scene text in the image?



The scene text is "THE EARTH WITHOUT ART IS JUST 'EH'".



What historical event is shown in the picture?



The picture shows the famous 1927 Solvay Conference, which was a landmark meeting in the history of physics. It was attended by many of the most prominent physicists of the 20th century, including Albert Einstein, Niels Bohr, Werner Heisenberg, Erwin Schrödinger, and Marie Curie. The conference was held in Brussels, Belgium, and focused on the development of quantum mechanics.

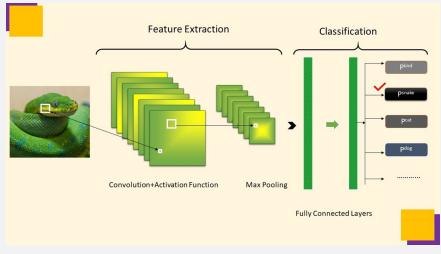
Figure 10: Qualitative examples of multimodal comprehension by SEED-X-I.

 Can this be true in the real world?



No, this is not possible in the real world. The image shows a stack of dice that are not physically connected, yet they are arranged in a way that suggests they are connected. This is an optical illusion, and the dice are not actually connected.

 Please explain the method in details.

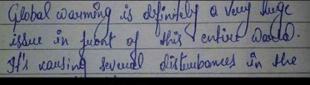


The method in the image is a Convolutional Neural Network (CNN) used for image classification. The process starts with the input image on the left, which is a green snake. This input goes through a convolutional layer, indicated by the green arrow pointing downwards. This convolutional layer is followed by a max pooling layer. The output of the max pooling layer feeds into a fully connected layer. This fully connected layer leads to the output of the network, which is the classification of the input image. The final output of the classification step is the classification of the input image, which is a green snake.

 What is the meaning of the meme?

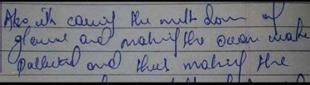
My Handwriting In Exams

First two pages



Global warming is definitely a very huge issue in front of this entire world. It's causing several disturbances in the

Middle pages



After carrying the math down of grown and making the ocean water polluted and that makes the

Last two pages



Yes, the black face mask located at green bbox can help protect you from catching the flu virus when you go out.

The meme is comparing the quality of handwriting in the first two pages of an exam and the last two pages. It suggests that the quality of handwriting often deteriorates as the exam progresses. This is a humorous way to express the common experience of writing neatly at the beginning of an exam and becoming less neat as the time progresses.

 Is there anything in the image that can protect me from catching the flu virus when I go out?



Yes, the black face mask located at green bbox can help protect you from catching the flu virus when you go out.

Figure 11: Qualitative examples of multimodal comprehension by SEED-X-I.