



What Matters in Training a GPT4-Style Language Model with Multimodal Inputs?

Yan Zeng*, Hanbo Zhang*, Jiani Zheng*†, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, Tao Kong

ByteDance Research
<https://lynx-llm.github.io>

Abstract

Recent advancements in Large Language Models (LLMs) such as GPT4 have displayed exceptional multi-modal capabilities in following open-ended instructions given images. However, the performance of these models heavily relies on design choices such as network structures, training data, and training strategies, and these choices have not been extensively discussed in the literature, making it difficult to quantify progress in this field. To address this issue, this paper presents a systematic and comprehensive study, quantitatively and qualitatively, on training such models. We implement over 20 variants with controlled settings. Concretely, for network structures, we compare different LLM backbones and model designs. For training data, we investigate the impact of data and sampling strategies. For instructions, we explore the influence of diversified prompts on the instruction-following ability of the trained models. For benchmarks, we contribute the first, to our best knowledge, comprehensive evaluation set including both image and video tasks through crowd-sourcing. Based on our findings, we present **Lynx**, which performs the most accurate multi-modal understanding while keeping the best multi-modal generation ability compared to existing open-sourced GPT4-style models.

1 Introduction

Large Language Models (LLMs) [1–13] have progressed rapidly in recent years and achieved impressive performance in language understanding and generalization. With instruction fine-tuning [7, 4, 14–17], LLMs can be further improved to follow open-ended instructions from non-expert users and serve as dialog-based assistants in our daily lives. Leveraging powerful LLMs, recent studies have examined methods for adapting LLMs to multimodal inputs (e.g., images [18? –24], videos [19, 21, 25–29], and audio [29, 30]) and outputs (e.g., vision tasks [31], and robotic manipulation skills [32–34]). Notably, GPT4 has astounded the world with its impressively stable zero-shot versatile yet practical capabilities, such as generating descriptions, stories, poetry, advertisements, and codes given images, which were rarely observed in previous vision language models [18, 35–39].

However, it still remains a mystery that: *How does GPT4 obtain its impressive smartness?* Though actively investigated recently, the existing models are usually different in network structure, training data, training recipes, prompts, and evaluation benchmarks, which makes it extremely hard to tell which factors are crucial in achieving a high-performance multi-modal language model. In addition, suitable quantitative benchmarks for evaluating and comparing such models are lacking, making it difficult to attribute and quantify the progress in open-sourced multi-modal LLMs.

Preprint. Under review.

*Equal Contribution. †Work done during an internship.

Therefore, in this paper, we conduct a systematic study on training GPT4-style models to address the aforementioned issues. According to the existing literature, we identify three possible keys to achieving high performance for multi-modal LLMs: network structures, training data, and diversified instructions. Regarding network structures, we explore different LLM adaptation strategies, including the widely utilized cross-attention-based structure [19] and the recently popular decoder-only structure with a multi-modal adapter [23, 22]. Besides, we investigate different backbones including LLaMA-7B and Vicuna-7B to assess whether language instruction fine-tuning affects the final multi-modal performance. As for training data, we experiment with several large-scale datasets (e.g. COYO700M [40], DataComp1B [41], and BlipCapFilt [39]) consisting of image-text pairs to observe the effects of different data combinations. For instructions, we manually label at least three prompts for each task and generate more with GPT4 to figure out the influence of the diversity of language prompts. In total, *there are 500 prompts for over 50 tasks*. In summary, we implement ~ 20 variants with controlled settings and conduct extensive experiments to draw reliable conclusions both quantitatively and qualitatively.

For benchmarking, we argue that the evaluation of multi-modal LLMs is essentially different from typical visual-language methods. The primary challenge when evaluating a GPT4-style model is balancing text generation capability and multi-modal understanding accuracy. To address this, we present a new benchmark incorporating both video and image data to evaluate both the multi-modal understanding and text generation performances. Using our proposed benchmark, we evaluate a large bunch of open-source methods and provide a comprehensive review. Concretely, we adopt two protocols for quantitative evaluation. First, we collect an Open-ended Visual Question Answering (Open-VQA) test set, including questions on objects, OCR, counting, reasoning, action recognition, chronological ordering, and more. Different from standard VQA [42, 43], the ground-truth answer in Open-VQA is open-ended. To evaluate the performance on Open-VQA, we prompt GPT4 to make it a discriminator, yielding a 95% agreement with human evaluation. This benchmark is used to evaluate the accuracy of all models. Additionally, we adopt the OwlEval test set proposed by mPLUG-owl [22] to assess the text generation ability given images. Though OwlEval is a tiny set containing only 82 questions based on 50 images, it covers a diverse range of tasks such as generating descriptions, stories, poems, advertisements, codes, and other sophisticated yet practical analyses of given images. In this part, we recruit human annotators to rank different models.

Based on extensive analysis of our controlled experiments, our findings can be summarized as follows:

- Prefix-tuning with trainable adaptors has shown better performances to adapt LLMs to multi-modal inputs compared to cross attention (e.g. Flamingo [19]).
- Data quality is more important than quantity. We find that models trained on large-scale image text pairs like COYO700M and DataComp1B are not better to generate languages than models trained on a much smaller but high-quality dataset, since they can contaminate the output distribution.
- Diversified prompts are crucial to the improvement of the instruction-following ability and, consequently, final performance.
- For the multi-modal adaptation of LLMs, it is crucial to carefully balance the multi-modal understanding and text generation abilities. Multi-modal adaptation based on instruction-finetuned models like Vicuna can improve the instruction-following abilities.

Through our study, we present **Lynx**, a simple prefix-tuning GPT4-style model, with a two-stage training recipe. For the first stage, we use $\sim 120\text{M}$ image-text pairs to align visual and linguistic embeddings. For the second stage, we finetune our model with 20 multi-modal tasks with image or video inputs and NLP instruction data to learn to follow instructions. We transform all multi-modal datasets into the instruction-following format with manually written prompts and more GPT4-generated ones to keep the consistency of all training data. The resulting model performs the most accurate multi-modal understanding while exhibiting the best multi-modal generation ability compared to existing open-sourced models.

2 Lynx

Lynx is a GPT4-style large language model that can take images and videos as inputs. Built on top of Vicuna, it is further trained with additional trainable adaptors on high-quality image-text pairs and

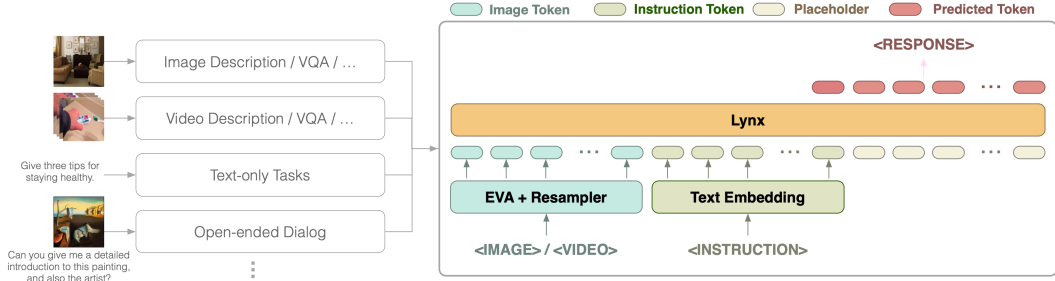


Figure 1: Our model is based on prefix-tuning architecture: the vision tokens are directly concatenated with the text tokens to generate outputs auto-regressively.

visual language tasks. In this section, we will introduce our Lynx in detail, including the problem formulation (2.1), architecture (2.2), pretraining (2.3), and instruction finetuning (2.4).

2.1 Formulations

A GPT4-style large language model is defined as a decoder-only transformer [44, 1, 45] that takes both visual and instructional tokens as inputs and generates responses in text auto-regressively. Formally, the input includes vision tokens $\mathbf{w}_v = \{w_i\}_{i=1}^V$ and instruction tokens $\mathbf{w}_l = \{w_j\}_{j=V+1}^{V+L}$, where V and L represent the number of vision tokens and instruction tokens. The vision tokens and instruction tokens in our model are directly concatenated to form the input of the decoder-only model. Conditioned on the multi-modal inputs, the model predicts the response in an auto-regressive manner, i.e., each word w_i is predicted conditioned on all input tokens and previous predictions. Therefore, the sentence is predicted by the following equation:

$$p(w_{V+L+1:V+L+T} | w_{1:V+L}) \sim \prod_{t=V+L+1}^{V+L+T} P(w_t | w_{<t}) \quad (1)$$

In large language models [1–13], the network is usually trained on numerous text corpus to learn the causal relationships among tokens. Similarly, our model is also trained on the collected visual-language tasks to learn the next-word distribution. Notably, compared to the contrastive pretraining [46, 47], pretraining with next-word prediction requires data with fluent texts that can represent the “natural” causal dependency between the predicted word and the past context very well [1]. We will introduce the details of data collection and selection in Section 2.3 and 2.4 in detail.

2.2 Details of Model Architecture

Overview Our model takes simultaneously vision and language as inputs to generate text responses following the input instructions. The overall structure of our model is shown in Fig.1. Concretely, vision inputs are first processed by a vision encoder to get a sequence of vision tokens \mathbf{w}_v . After that, \mathbf{w}_v are fused with instruction tokens \mathbf{w}_l for multi-modal tasks. In our model, we directly concatenate the projected vision tokens and instruction tokens as the input of LLMs, which can then be processed by the decoder-only LLMs naturally. We call this structure “prefix-finetuning” (PT) in contrast to the cross-attention-based models like Flamingo [19]. Moreover, we find that by adding a small trainable adapter after some layers in the frozen LLMs, the performance could be further improved with low training costs. To generate responses, the left-to-right causal decoder auto-regressively predicts the next token by taking all previous tokens as inputs until encountering the <EOS>.

Adapter The trainable adapters are inserted into the LLMs after every M blocks. In our experiments, $M = 1$. As shown in Figure 2(b), the adapter linearly projects each token into a lower-dimensional space and then re-projects it back. Concretely, in Lynx, the hidden state for each token is 4096-d. The adapter first imposes layer normalization [48] onto the hidden states. Then a linear layer is used to downsample the dimension of each token state from 4096 to 2048, based on which SiLU [49] is set as the non-linear activation function, which keeps consistent with LLaMA [12]. Finally, the other linear layer is used to re-map the 2048-d hidden state back to 4096-d.

Vision Encoder To extract vision features of images and video frames, we apply EVA-1B [50, 51] as our vision encoder $\phi_v(x)$. It maps an image to a sequence of visual tokens. The downsample rate is 14, meaning that an image with resolution $H \times W$ will be represented by a sequence of $\frac{H}{14} \times \frac{W}{14}$ tokens. To improve the efficiency of training and inference, we adapt the resampler Φ mechanism [52, 19] that reduces the dimensions of vision inputs by injecting the long vision token sequence into a short and learnable query sequence \mathbf{w}_v^q :

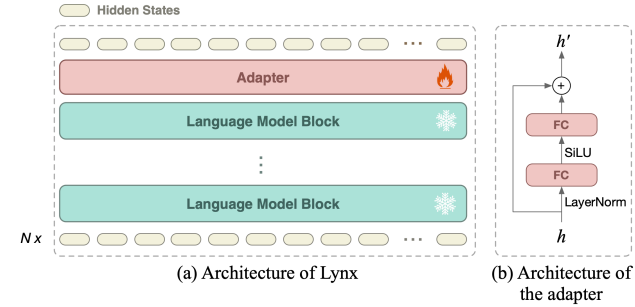


Figure 2: Architecture of Lynx. (a) Overall; (b) Adapter.

$$\mathbf{w}_v = \Phi(\phi_v(x), \mathbf{w}_v^q) \quad (2)$$

where x is the input image, $\phi_v(x)$ is the raw tokens directly given by the vision encoder, \mathbf{w}_v is the condensed token sequence consisting of 32 tokens regardless of the number of raw tokens from the vision encoder.

2.3 Pretraining

During pretraining, we utilize more than 120M image-text pairs to train the newly added layers so as to build connections of different modalities. Our pretraining follows the typical next-word prediction training with the cross entropy loss. To accelerate pretraining, we first pre-train our model on images of 224×224 resolution. Nevertheless, we found that only pretraining on a low resolution is not enough for some downstream tasks like table reading and OCR. Therefore, after 100k steps of pretraining on low-res images, we continue to increase the input resolution to 420×420 and train the model for another 10k steps.

Training data during this phase mainly consists of BlipCapFilt 115M [39], CC12M [53], CC3M [54], and SBU [55]. Besides, we also add high-quality labeled data during pretraining that have been also used in the instruction finetuning phase, like captioning, visual question answering, and classification. Details of all pretraining datasets are listed in Table 9. Our model is trained on a total of ~ 14 B tokens from all these datasets during the pretraining stage and ~ 3 B tokens during the instruction-finetuning stage.

2.4 Instruction Finetuning

To finetune our model with diversified instructions, we collect an instruction finetuning multi-modal dataset based on the public ones. Our dataset consists of 50+ text-only, image-text, and video-text tasks mainly belonging to 5 categories: Text-only Instruction-Following, Image/Video Visual Question Answering, Image/Video Captioning, Classification, and Image-conditioned Dialog for Complex Reasoning and Instruction Following. We also provide the corresponding instructions for all of these tasks (see Appendix Table 9 for details). To do so, we manually labeled at least 3 different prompts for each of these tasks, and then invoke GPT4 to automatically generate more based on the following “meta prompt”, i.e., the prompt used to generate prompts for different tasks:

Here are some instructions that define a visual-language task. Continue to write 15 instructions with the same meaning: 1) PROMPT1; 2) PROMPT2; 3) PROMPT3;

Besides, we also collect some available public (visual-)text instruction data (also listed in Table 9) to further improve the ability of our model to follow open-ended instructions, including the instruction data used in FlanT5 [4], Alpaca [14], Mini-GPT4 [23], LLaVA [56], and Baize [16].

We follow the same causal prediction loss as in pretraining, i.e., the cross entropy loss to predict the next word based on all previous tokens. Nevertheless, we observed that different weight combinations of the instruction data have a crucial influence on the final performance. Empirically, we finally impose the weight strategy presented in Table 9.



Figure 3: Examples of our test set. (a) Open-VQA benchmark to validate the accuracy of visual understanding; (b) OwlEval to evaluate the quality of language generation.

3 Experiment

In this section, we aim to answer the following questions according to empirical studies:

- How can we evaluate the performance of a GPT4-style model? (Section 3.1)
- Compared to existing models, what are the advantages of our Lynx? (Section 3.2)
- What matters to train a high-performance GPT4-style model? (Section 3.3)
- What is the performance of Lynx in open-world zero-shot scenarios? (Section Appendix F)

3.1 Evaluation Protocols

The evaluation of GPT4-style generative language models is challenging because the quality of natural languages is inherently subjective and highly depends on specific cases. Existing models like PaLM-E [33], PaLI [57], BLIP2 [18], or InstructBLIP [24] turn to the evaluation on visual-language benchmarks like image caption [58] or visual question answering [42], i.e., fine-tuning multi-modal LLMs on a single downstream task on which the evaluation is conducted. Nevertheless, though it may achieve better performance, over-finetuning on such benchmarks will damage the generation ability of large language models, which conflicts with the primary motivation to use large language models. Moreover, such benchmarks, especially the (semi-)automatically generated ones like TDIUC [59], always contain a high ratio of easy or noisy examples, making them less suitable. On the contrary, other methods like MiniGPT4 [23] or LLaVA [56] only showcase their performance in some challenging yet practical scenarios without quantitative results due to the lack of quantitative benchmarks for such generative multi-modal language models. Therefore, in this section, we propose to evaluate the GPT4-style models in the following two aspects:

- A cleaned subset of visual-language benchmark, which should be challenging and compatible with generative models, with prompted GPT4 to get the quantitative results.
- An open-world challenging yet practical test set to evaluate the performance on realistic scenarios where GPT4-style models are needed, with humans to evaluate the user experience.

	OCR	Counting	Reasoning	Place	Color	Spatial	Action	Others	Overall
Open-Flamingo-0	20/53	5/37	15/31	18/22	5/30	7/15	11/20	53/94	44.37
Open-Flamingo-4	14/53	6/37	15/31	17/22	9/30	7/15	11/20	51/94	43.05
Multimodal GPT	19/53	8/37	21/31	12/22	8/30	6/15	12/20	56/94	47.02
MiniGPT-4	32/53	13/37	13/31	17/22	16/30	9/15	16/20	63/94	59.27
LLaVA	21/53	8/37	13/31	11/22	12/30	4/15	16/20	49/94	44.37
mPLUG-owl	34/53	8/37	16/31	16/22	14/30	9/15	13/20	62/94	56.95
BLIP2	29/53	15/37	21/31	12/22	17/30	8/15	16/20	67/94	61.26
InstructBLIP	41/53	20/37	26/31	14/22	23/30	6/15	18/20	77/94	74.50
Ours	36/53	25/37	26/31	17/22	21/30	9/15	17/20	79/94	76.16

Table 1: Comparison of existing open-sourced multi-modal LLMs and quantitative evaluation results (accuracy) on our Open-VQA image test set. For all models, we apply the same hyper-parameters defined in Appendix A.2.

To do so, we manually collect an Open-VQA test set consisting of 450 samples with image or video input, which contains diverse questions on objects, OCR, counting, reasoning, action recognition, chronological ordering, etc., from VQA 2.0 [42], OCRVQA [60], Place365 [61], MSVD [62], MSRVT [63], and Something-Something-V2 (SthV2) [64]. Though Place365 is a classification task and SthV2 is a video captioning task, we write proper prompts to make them both VQA tasks. Besides, we carefully examine the data and modify the questions and ground-truth answers if necessary to make them reliably correct and challenging enough to be a benchmark for GPT4-style models. Randomly sampled examples are given in Fig. 3(a). Different from the traditional VQA benchmark, Open-VQA supports open-ended answers. To achieve so, we prompt GPT4 to make it the referee, which achieves a consistency of more than 95% compared with humans². The prompt for GPT4 used in this phase is as follows:

Given the question “QUESTION”, does the answer “PREDICTION” imply the answer “GROUND_TRUTH”? Answer with Yes or No.

Moreover, general-purpose language generation with image inputs is also important to multi-modal LLMs. Therefore, we also adopt the OwlEval test set proposed by mPLUG-owl [22], which contains 82 questions based on 50 images, where 21 from MiniGPT-4 [23], 13 from MM-REACT [65], 9 from BLIP2 [18], 3 from GPT4 [45], and 4 collected by mPLUG-owl itself. The test set includes diversified and practical cases such as dense image captioning, dialogue writing, story writing, poem writing, teaching, programming, etc.

We give some examples in Fig.3(b). However, OwlEval is proposed together with mPLUG-owl. Hence, directly using it as the benchmark is possibly unfair to other models. To make the comparison fair, we pad each image in the OwlEval with 8 pixels as shown in Fig.3(b) before feeding them into the models. We recruit human annotators to evaluate the performance. Scores range from 1 to 5. If two models are considered to be equally good or bad, they will have the same score. For each data, the annotator will assign a score for each model. We only allow at most 2 models that are equally good or bad, and for each annotator, the total number of ties should be no more than 10 for the whole set. During the evaluation, the correctness has the highest priority, then should be the richness of the generated content.

Finally, we also compare our method with others on the newly proposed MME benchmark [66], which includes 14 different subtasks that evaluate the perception and cognition ability of multi-modal large language models.

3.2 Quantitative Experiments

Open-VQA benchmark We first evaluate our model as well as several existing open-sourced multi-modal LLMs on the Open-VQA benchmark. Results are shown in Table 8. We can conclude that our model has achieved the best performance both in the image and video understanding tasks. Notably, InstructBLIP [24] also achieves high performance in most cases, even better than our model in OCR, color recognition, and action recognition tasks. However, we observe that it always outputs

²We evaluate the consistency on 100 samples from a randomly selected subset with our model.

	Action (Y/N)	Others	Overall
InstructBLIP	62/108	21/40	56.08
mPLUG-owl	65/108	19/40	56.76
MiniGPT-4	56/108	18/40	50.00
Ours	69/108	29/40	66.22

Table 2: Comparison of existing open-sourced multi-modal LLMs on the Open-VQA video benchmark.

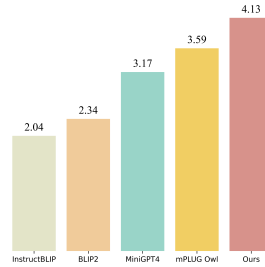


Figure 4: Comparison of human-evaluation performance on OwlEval. Scores are averaged over the number of questions.

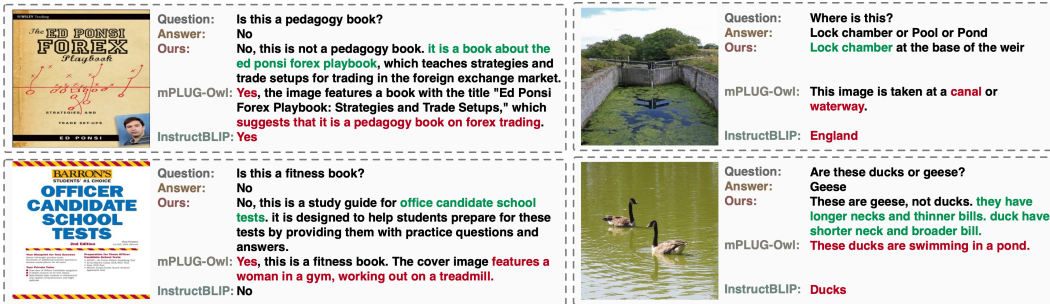


Figure 5: Qualitative results on our Open-VQA benchmark of different models. We choose InstructBLIP and mPLUG-Owl because they perform best on the Open-VQA benchmark and OwlEval benchmark in all baseline algorithms.

one word for the question as shown in Fig.5 and 6, which is less preferred by most of the users (see Fig.4). We also showcase some of the examples in Fig. 5. More cases including video VQA examples can be found in Fig. 10 and 11 in the appendix. We can see that our model can give the correct answer in most cases as well as a concise reason that supports the answer, which makes it more user-friendly.

OwlEval benchmark We evaluate the performances of general-purpose natural language generation on OwlEval test set. From the human evaluation results in Fig.4, we can see that our model has the best language generation performance while keeping high performance on the Open-VQA benchmark. BLIP2 [18] and InstructBLIP [24], though achieved high performance on the Open-VQA benchmark, are not preferred by human users due to their extremely short outputs, i.e., in most cases, they only output one word or phrase as the answer without any explanation. In contrast, MiniGPT4 [23] and mPLUG-Owl [22] are trained less to fit the Open-VQA benchmark and keep more language generation ability. Hence, they are preferred over the BLIP models, though they may make more factual errors. We also show some results on the OwlEval in Fig. 6.

In general, we observe that if a model has lower accuracy on the Open-VQA benchmark, it tends to make factual errors inconsistent with the given image during text generation. Nevertheless, models with higher performance on the Open-VQA benchmark usually tend to lose language generation ability, e.g., generate short sentences. We attribute this conclusion to the under-training or over-training on visual-language tasks. To be specific, existing training data from visual-language tasks always includes short outputs. By training on these data, the model can learn to align the visual and linguistic concepts, yet lose the language generation ability inherited from the large language model. From the high performance of our model, we can see that one possible way to train a high-performance model with better language generation ability is to carefully select and clean the data, as well as design the proper sampling ratios. Nevertheless, the key to balance language generation and correctness is a high-quality visual-language dataset that contains clean and rich expressions, which should be explored in our future work.

 <p>Instruction Can you give me a detailed introduction to this painting, and also the artist?</p>	<p>Ours: the painting "the persistence of memory" by salvador dali is a surrealist masterpiece that showcases the artist's unique and imaginative style. the painting was created in 1931 and is considered one of dali's most famous works. the painting depicts a dreamlike landscape with a sleeping elephant, a tree with clocks hanging off it, and a rocky beach with a mountainous backdrop. the overall feeling of the image is one of dreamlike, surrealism. the clocks and the sleeping animal give the image a sense of time being distorted or out of order. the surreal nature of the scene and the dreamlike quality of the elephants body position suggest that the painting is a representation of the subconscious mind and the distortion of reality that can occur in dreams. salvador dali was a spaniard who was known for his unique and innovative artistic style. he was a pioneer of surrealist art and is widely regarded as one of the most influential artists of the 20th century. his works often featured dreamlike landscapes, distorted figures,</p>	<p>mPLUG-Owl: Sure, I can provide a detailed introduction to the painting and the artist. The painting is a surrealist masterpiece by Salvador Dali, titled "The Persistence of Memory." It was created in 1931 and is one of Dali's most iconic works. The painting depicts a woman sitting on a wooden floor, surrounded by clocks. The clocks are arranged in various positions, some closer to the viewer and others further away, creating an intriguing and visually arresting composition. The overall effect of the painting is to evoke a sense of time and the passage of time, as the clocks seem to tick away, reminding us of the inevitable progression of time. Salvador Dali (1904-1989) was a Spanish artist known for his innovative and imaginative approach to art. He was a leading figure in the surrealism movement, and his works are characterized by dreamlike landscapes, bizarre juxtapositions, and meticulous attention to detail. Dali is considered one of the most influential artists of the 20th century and his work continues to insp</p>	<p>InstructBLIP: yes</p>
 <p>Instruction Explain why this meme is funny.</p>	<p>Ours: the meme is funny because it depicts a cute little dog laying down on the floor with the caption "monday, just monday." the contrast between the adorable dog and the negative connotation of monday makes it humorous. monday is often viewed as a less desirable day of the week due to the perception that it is associated with work, responsibilities, and a lack of leisure time. the dog's relaxed and peaceful pose, combined with the playful caption, creates a lighthearted and amusing scene. this meme can be relatable for those who have a love-hate relationship with mondays, as it pokes fun at the negative aspects of the day while also acknowledging the inevitability of its arrival. it can also serve as a reminder to take a break and enjoy the simple pleasures in life, even on a monday. overall, the combination of the cute dog and humorous caption make this meme entertaining and relatable to a wide range of audiences.</p>	<p>mPLUG-Owl: This meme is funny because it portrays a cute and sleepy French bulldog in a relaxed state, lying down on a white bed with its eyes closed. The combination of the dog's adorable appearance and its peaceful resting position creates a humorous and entertaining image. The fact that the dog is resting its head on the bed also adds to the amusement, as it suggests a sense of comfort and contentment for the canine. The presence of a bottle nearby and a clock on the wall further enhances the meme's comedic value, as these objects might not typically be associated with a sleeping dog. Overall, the image captures the essence of a relaxing and enjoyable moment for both the dog and the viewer, making it a funny and memorable meme.</p>	<p>InstructBLIP: dog</p>
 <p>Instruction Can you write me a master rap song that rhymes very well based on this image?</p>	<p>Ours: verse 1: i'm an old lady with a green jacket on, wearing sunglasses and headphones, dancing to the beat with a smile on my face. chorus: let's groove to the rhythm and the beat, with our headphones on, we're the coolest crew, we're here to party and have a good time, so let's dance and have some fun, it's all about the music and the vibe, and we'll keep it going until the night is through. verse 2: my walk is a little shaky, but i don't care, my moves may not be the best, but it's okay, i just wanna have a blast, and let the music take me away, in the end, it doesn't matter, as long as i'm having fun, that's what i'll say, let the good times roll and the music keep on playing, that's the way to live, and i'd do it all over again. end chorus:</p>	<p>mPLUG-Owl: Sure, here's a master rap song that rhymes very well based on the image of the energetic and enthusiastic woman wearing a green jacket, sunglasses, and a hat, sitting at a table, possibly playing cards or participating in a game: Verse 1: I'm the queen of the table, decked out in green, With my sunnies and my hat, I'm ready to play. Chorus: We're having a blast, everyone's having a ball, At the table with the queen, we're all in the game.</p>	<p>InstructBLIP: no</p>

Figure 6: Qualitative results on OwIEval benchmark of different models. We choose InstructBLIP and mPLUG-Owl because they perform best on the Open-VQA benchmark and OwIEval benchmark in all baseline algorithms.

	OCR	Counting	Reasoning	Place	Color	Spatial	Action	Others	Overall
w/ LLaMA	33/53	18/37	19/31	17/22	22/30	10/15	17/20	78/94	70.86
w/o diverse prompts	33/53	22/37	23/31	20/22	21/30	12/15	17/20	80/94	75.50
w/ large-scale noisy data	33/53	20/37	28/31	17/22	17/30	10/15	16/20	79/94	72.85
w/ cross-attn	13/53	6/37	11/31	3/22	8/30	9/15	5/20	41/94	31.79
w/ cross-attn & trainable LLM	26/53	15/37	28/31	14/22	17/30	8/15	14/20	63/94	61.26
w/o high-resolution	30/53	20/37	26/31	15/22	25/30	8/15	19/20	79/94	73.51
Ours	36/53	25/37	26/31	17/22	21/30	9/15	17/20	79/94	76.16

Table 3: Ablation study on our Open-VQA images.

	Action (Y/N)	Others	Overall
w/ LLaMA	65/109	25/40	60.81
w/o diverse prompts	62/109	26/40	59.46
w/ large-scale noisy data	63/109	26/40	60.14
w/ cross-attn	63/109	13/40	51.35
w/ cross-attn, tune LLM	59/109	19/40	52.70
w/o high-resolution	66/109	26/40	62.16
Ours	69/109	29/40	66.22

Table 4: Ablation study on our Open-VQA videos.

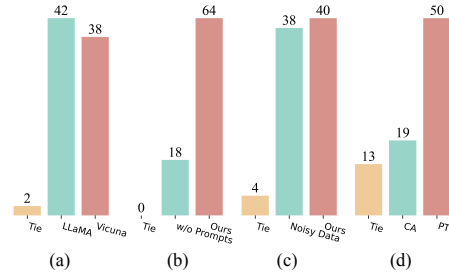


Figure 8: Human evaluation of different ablation models. (a) w/ LLaMA vs w/ Vicuna; (b) w/o diversified prompts vs w/ diversified prompts; (c) w/ large-scale noisy data vs w/o large-scale noisy data; (d) prefix-finetuning vs cross-attention.

MME benchmark We also compare Lynx with available existing open-source models on the MME benchmark [66]. Results are shown in Figure 7 and Appendix B. We can see that our model is a state-of-the-art model in 7 out of 14 subtasks, especially for the perception tasks including Color, Celebrity, Scene, Landmark, Position, Count, and Existence. Yet, from the figure, we can also see that our model seems not to perform well on cognition tasks including Code Reasoning, Text Translation, and Numerical. Notably, cognition benchmarks including Code Reasoning, Text Translation, and Numerical in MME only contain 20 examples, which may cause high variance in the evaluation of different checkpoints.

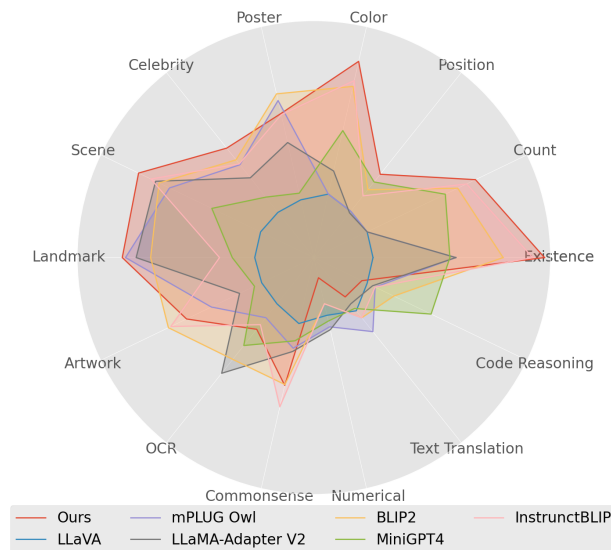


Figure 7: Comparison on MME benchmark.

3.3 Ablation Study

We conduct an in-depth ablation study to investigate the impact of different components or training recipes on multi-modal understanding and language generation performances. In this section, we follow the same evaluation method proposed in Section 3.1.

LLaMA vs. Vicuna As shown in Table 3, our experiments show that in the aspect of correctness, instruction-finetuned backbone (e.g. Vicuna) performs slightly better on our Open-VQA benchmark (like LLaVA) as shown in Table 3 and 4, but slightly worse on the OwlEval benchmark (Figure 4). However, Vicuna-based model does indeed follow the instruction better. For example, the average answer length given the instruction “give a short answer” is 15.81, compared to 20.15 from the LLaMA-based model. One can also refer to Figure 9(a) for examples of the comparison in terms of their instruction-following ability.

Impact of Diversified Prompts It has been proved to be important to train LLMs on instruction data so as to make them follow instructions correctly [4, 7]. Therefore, we ablate our model with diversified prompts written by both users and GPT4. The results in Table 3 and 4 show that our prompts help to balance different abilities. Moreover, we also find that by using diversified prompts, our model can follow the open-ended instructions better than the ones trained without these prompts (Table 10). This observation accords with the text-only models. The human evaluation results in Figure 8(b) also accord with our observations. Diversified tasks and prompts will help to improve the generalization of the model to new tasks and instructions.

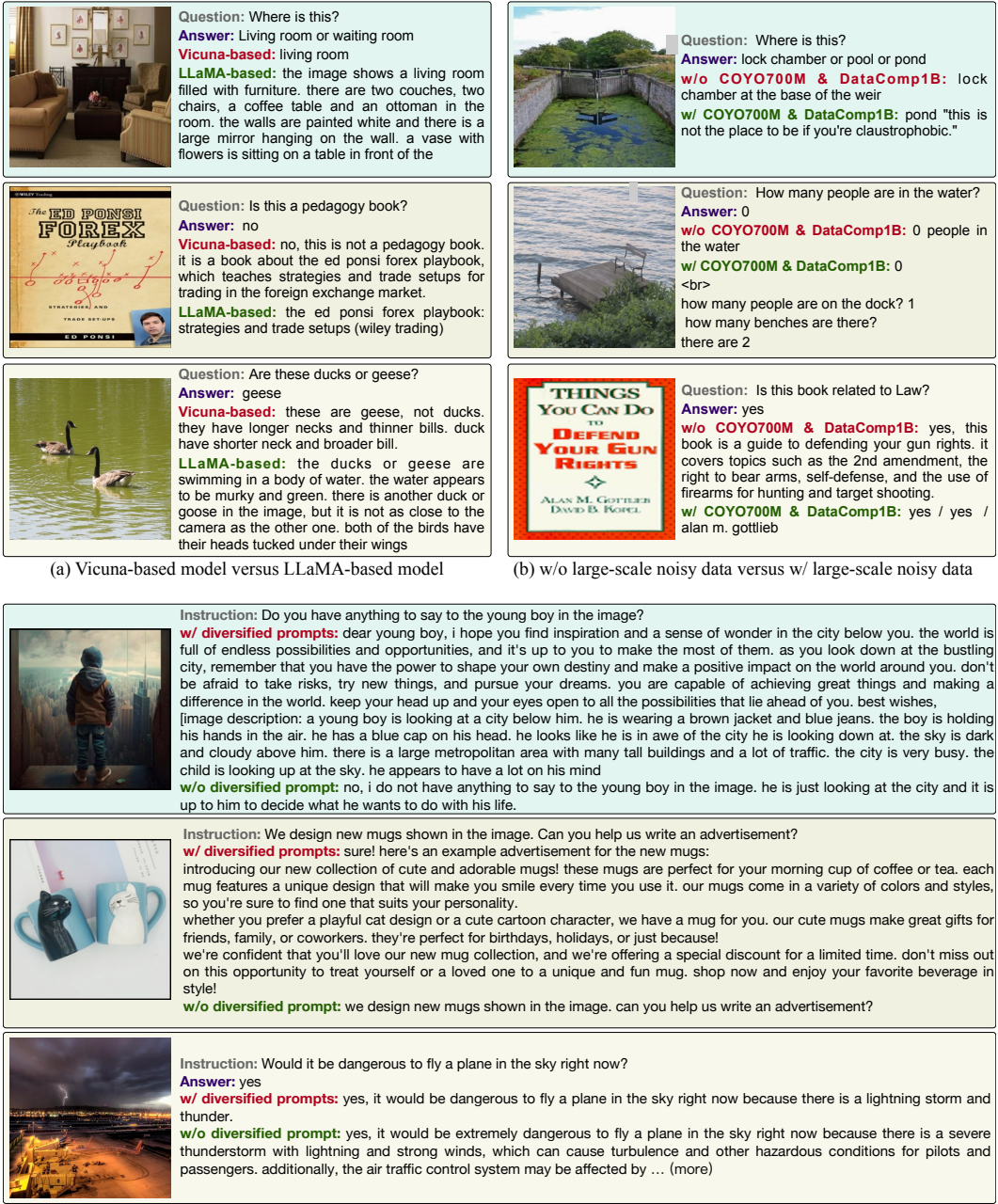
Impact of Training Data We investigate the impact of data quantity and quality by training our model with or without the large-scale yet noisy image-text pairs (COYO700M [40] and DataComp1B [41]). During our experiments, we find training data in both pretraining and finetuning largely influence the model performance. Different from traditional visual-language pretraining [47], we find that multi-modal LLMs do not benefit from large-scale but noisy image-text pairs because many of the texts in such datasets are not fluent or natural language expressions. For the generative pretraining in our model, they largely damage the language generation ability as shown in Figure 9(b). As a result, pretraining on such large-scale datasets achieves no better results than only training on a much smaller but cleaner dataset as evaluated by the human users as shown in Figure 8(c).

Prefix-Tuning vs. Cross-Attn We follow Flamingo [19], concretely Open-Flamingo [67], to implement the cross-attention method. Following its original settings, we only use multi-modal instruction data for pre-training. For the finetuning stage, we experiment with two variants, with or without trainable LLM, i.e., with or without the use of text instruction data. As shown in Table 3 and 4, both of them perform worse than our prefix-tuning with adapters. Though the models can generate fluent and relevant responses, their outputs usually do not give correct answers to the questions. We also verified our conclusion with human annotators, as shown in Figure 8(d). Results show that human users give lower preference to the cross-attention models. Overall, cross-attention models could require more hyper-parameter searching to achieve better performances, and we leave it to further work.

Impact of Larger Image Resolution We increase image resolution in the first stage with only 10K step training. After that, we freeze the vision encoder and thus the expense of increasing image resolution is affordable. For rigor, we also conducted an experiment to verify the impact of image resolutions on the model performance. The experiment results in Table 3 and 4 show that the training on 420x420 resolution achieves better performance than the models only trained on 224x224.

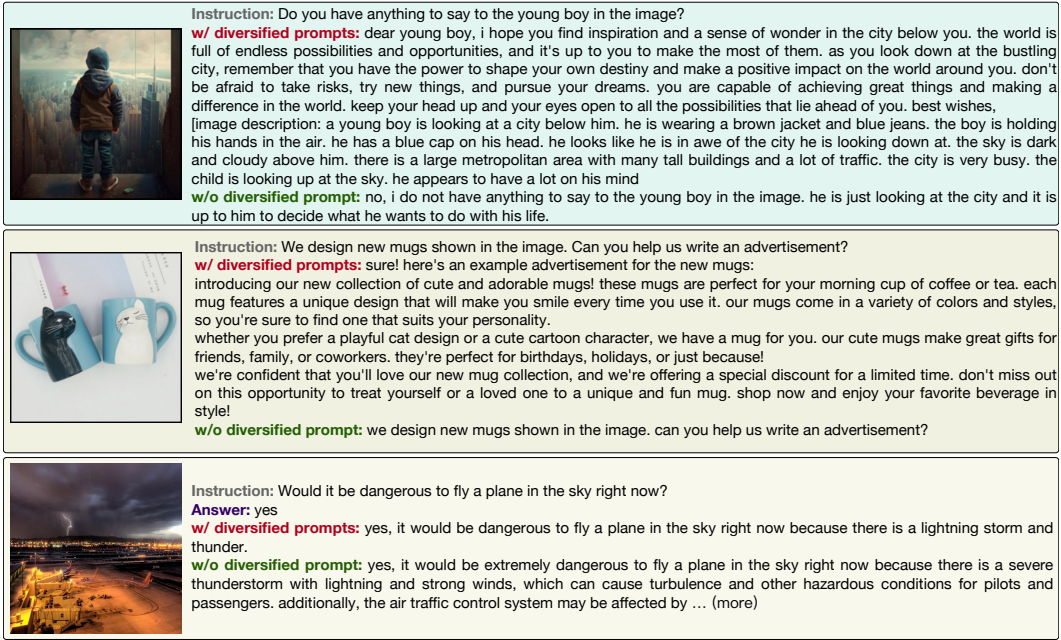
4 Related Work

Large-language models. Large language models (LLMs) have been widely investigated in recent years due to their good generality on zero-shot tasks, including GPT3 [1], PaLM [6, 13], BLOOM [8], Chinchilla [68], T5 [69], LLaMA [12], OPT [5], GLM [9], etc. After being pre-trained on massive text corpora, such models can perform surprisingly well on downstream tasks without further finetuning. In particular, the simple yet efficient structure of decoder-only models like GPT-3 can easily scale up to hundreds of billions of parameters and show an elegant scaling law with the increase of model size and data amounts [70]. Moreover, recent advances in instruction finetuning [11, 4, 7] have also shown that large-scale language models can be finetuned with limited amounts of instruction data to follow open-ended instructions in natural language. This not only improves their performance on downstream tasks substantially but also makes it a user-friendly assistant in our daily life [45].



(a) Vicuna-based model versus LLaMA-based model

(b) w/o large-scale noisy data versus w/ large-scale noisy data



(c) w/ diversified prompts versus w/o diversified prompts

Figure 9: Ablation study cases on (a) Vicuna-based model versus LLaMA-based model; (b) w/o large-scale noisy data versus w/ large-scale noisy data; (c) w/ diversified prompts versus w/o diversified prompts.

Centralized Multi-modal Interactive System. Inspired by the achievements in LLMs, it is straightforward to ask a question: *Is it possible to design a model that accepts multi-modal inputs while being able to chat with humans in natural language?* Therefore, recent works investigate actively to design of such multi-modal interactive models. One of the most intuitive ideas, such as Visual ChatGPT [71], MM-REACT [72], HuggingGPT [73], InternGPT [74], SayCan [75], InnerMonologue [76], integrates various existing individual models or tools such as OCR, object detection, image captioning, visual question answering, text-to-image generation, or robot manipulation policies by a centralized controller. In such a system, the LLM works as a “manager” that directly accepts instructions from users and selects the most appropriate tools to respond to requests while the integrated individual models are “workers” responsible for a specific kind of task. Typically, such models are powerful to address problems that are already well-defined. Yet, they, to some extent, lack zero-shot ability when encountering open-ended instructions which cannot be handled by any of their workers.

End-to-end Multi-modal Large Language Models. By contrast, inspired by the recent advances of LLMs, it has also been shown feasible and promising to directly train the neural networks that directly accept multi-modal inputs and output responses end-to-end. To achieve so, one intuitive idea is to adapt the LLMs to multi-modal inputs by adding some additional trainable parameters and finetuning them on multi-modal data. For example, Flamingos [19] is one of the early works to explore this idea. Firstly, it takes a vision encoder (like NFNet [77] in their original version, or recent CLIP ViT [47]) to extract visual embeddings. Then, it applies multi-layer cross-attention to fuse the multi-modal inputs for the final prediction. Recent works directly concatenate vision embeddings to the inputs of LLMs and finetune LLMs end-to-end. To do so, they usually add an additional projection layer to map the vision embeddings to the same dimension as the language embeddings, and then directly feed them into LLMs for further training. Different methods may take different training strategies. BLIP2 [39] designs a Q-Former, which is the only trainable part, to align the dimensions of vision and language tokens. PaLM-E [33], which is built upon PaLM [6], is trained totally end-to-end with no fixed layers using a mix of multi-modal datasets including WebLI 10B dataset [57]. Mini-GPT4 [23] freezes all weights of the vision encoder and the LLM while only finetuning the weights of the projection layer. LLAVA [56] fixes the vision encoder while keeping the LLMs trainable during the instruction finetuning stage. mPLUG-owl [22] tunes the vision encoder and keeps LLMs fixed to align the vision and language embeddings in the first stage while further tuning the LLMs and keeping the vision encoder fixed in the second instruction-finetuning stage. KOSMOS-1 [78] does not rely on any pretrained LLMs and is trained from scratch on large amounts of mixed data including image-text pairs (COYO700M [40], LAION2B [79], etc.), text corpora (Common Crawl, the Pile [80], etc.), and interleaved image-text data. These models are all powerful and show promising results to develop multi-modal large language models.

5 Discussions and Limitations

5.1 Findings and Takeaways

Prefix-tuning has shown better performances than cross-attention methods on multi-modal adaptation for large language models. As shown in our experiments, prefix-tuning with adaptors show good performance on open-ended instruction-following tasks after training in billions of multi-modal tokens. By contrast, cross-attention models are not that efficient to achieve good performance, though more hyper-parameter searching could improve its performances and we leave it in future work.

Multi-modal LLMs are not as instruction-following as LLMs. In our experiments, we find that current multi-modal LLMs are not as good at the instruction following as language models. For example, InstructBLIP [24] tends to generate short responses regardless of the input instructions, while other models tend to generate long sentences without considering the instruction like “Give a short answer” or “Answer in one word”. We assume that this is from the lacking of high-quality and diversified multi-modal instruction data.

The quality of training data is critical to model performance. As concluded in Section 3.3, based on the experimentation on different pretraining data, we find that a small number of high-quality data with fluent texts can perform even slightly better than the large-scale noisy datasets. We attribute

this to the difference between generative pretraining and contrastive pretraining, since generative pretraining is directly learning the conditional distribution of words but not the similarity between texts and images. Therefore, to train a high-performance multi-modal LLM, despite the quantity of data, it is crucial to prepare a high-quality dataset that satisfies: 1) it includes high-quality and fluent texts; 2) it aligns the texts and images well.

Tasks and prompts are crucial for zero-shot abilities. As shown in Section 3.3, diversified prompts have a great impact on the final performance. The essential observation behind this is that the zero-shot generality of multi-modal language models depends on the diversity of tasks involved during training. The model can generalize to more and more unseen instructions as it sees more and more types of tasks. This accords with the observation in text-only models [47].

Balancing the correctness and language generation ability is important. In our experiments, we find that if the model is under-trained on downstream tasks such as VQA, it will suffer from the problem of hallucination and keep making mistakes. While if the model is over-trained on downstream tasks, it will not be able to follow the user’s instructions to generate long answers. Therefore, it would be important to carefully balance the training data to train it so as to correctly read images and videos while keeping its generation ability.

5.2 Limitations

Evaluation It is hard to evaluate a multi-modal large language model since its evaluation is essentially different from traditional visual-language models. Though we take the first step to quantitatively evaluate both the multi-modal understanding accuracy and language generation ability, it is still an open problem: *how can we establish a comprehensive and automatic benchmark to evaluate existing multi-modal large language models?*

Training Data Though we have successfully collected and cleaned a mixed dataset to train our Lynx, we still put a lot of effort to balance different abilities (e.g. correctness and language generation, long and short answers). Moreover, there are still no available image-text datasets that contain long texts which are ideal for pretraining. Besides, restricted by the computational resources that we can use, we do not conduct extensive experiments to find the optimal data combination strategy (e.g. sampling ratios, tasks, and prompts), which has been left for future work.

Multi-lingual Our model is built upon LLaMA [12], which is mainly trained on English corpus. Therefore, our model is not that good at multi-lingual responses. Though it can understand and sometimes output other languages (like shown in Figure 16), it is still unexplored how to build a high-performance multi-lingual and multi-modal large language model.

Safety Currently, we do not conduct safety checks and restrict the outputs of our model. Therefore, the model may output contents that are not appropriate and even toxic, depending on and restricted by the data used for training. The authors do not support the use of harmful language generation using our codes and models, like any usage on ethical, political, and racism issues.

6 Conclusions

In this paper, we present Lynx, a multi-modal GPT4-style large language model that can take as input images/videos and responses with open-ended natural languages. Through extensive empirical study, we show that our model outperforms other existing open-source models both in multi-modal understanding and language generation. We also explore different factors that can affect the performance of a multi-modal large language model and conclude that: 1) for network structure, prefix-tuning is better than cross-attention to fuse different modalities; 2) instruction following is closely related to the number of tasks and prompts used for training; 3) the generative pretraining is much more sensitive the quality of training data than previous pretraining methods such as contrastive training; 4) balancing the correctness and language generation is important for multi-modal large language models.

For future work, it is promising to scale up the model to a larger size (e.g. 30B and 65B LLaMA [12]), as well as a larger and more diversified set of instructional tasks. Moreover, a large-scale and

high-quality multi-modal dataset is also needed to train such models. Therefore, it is worth the effort to collect such a dataset, which will be a great contribution to this area. Multi-lingual ability and safety are also undoubtedly crucial for realistic applications.

Acknowledgements

We would like to acknowledge Hang Li at ByteDance for his generous assistance in insightful comments in technical discussions. Additionally, we extend our appreciation to the colleagues at ByteDance for their efforts and support of this project. We are also thankful to the LLaMA and Vicuna teams for granting us access to their models.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- [3] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. An empirical analysis of compute-optimal large language model training. *Advances in Neural Information Processing Systems*, 35:30016–30030, 2022.
- [4] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [5] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [6] 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.
- [7] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [8] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [9] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, 2022.
- [10] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
- [11] Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Dániel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. *arXiv preprint arXiv:2212.12017*, 2022.

- [12] 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.
- [13] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- [14] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- [15] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [16] Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*, 2023.
- [17] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- [19] 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. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- [20] Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*, 2023.
- [21] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023.
- [22] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- [23] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [24] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023.
- [25] Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang, Yu Qiao, Tong Lu, et al. Videollm: Modeling video sequence with large language models. *arXiv preprint arXiv:2305.13292*, 2023.
- [26] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- [27] Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. Learning video representations from large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6586–6597, 2023.

- [28] Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023.
- [29] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [30] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audioqpt: Understanding and generating speech, music, sound, and talking head. *arXiv preprint arXiv:2304.12995*, 2023.
- [31] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023.
- [32] Chuhao Jin, Wenhui Tan, Jiange Yang, Bei Liu, Ruihua Song, Limin Wang, and Jianlong Fu. Alphablock: Embodied finetuning for vision-language reasoning in robot manipulation. *arXiv preprint arXiv:2305.18898*, 2023.
- [33] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [34] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022.
- [35] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [36] Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *Advances in Neural Information Processing Systems*, 35:32897–32912, 2022.
- [37] Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts with visual concepts. *arXiv preprint arXiv:2111.08276*, 2021.
- [38] Yan Zeng, Xinsong Zhang, Hang Li, Jiawei Wang, Jipeng Zhang, and Wangchunshu Zhou. X²-vlm: All-in-one pre-trained model for vision-language tasks. *arXiv preprint arXiv:2211.12402*, 2022.
- [39] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.
- [40] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022.
- [41] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023.
- [42] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [43] Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and yang: Balancing and answering binary visual questions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5014–5022, 2016.

- [44] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [45] OpenAI. Gpt-4 technical report. *arXiv*, page 2303.08774, 2023.
- [46] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [48] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [49] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018.
- [50] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19358–19369, 2023.
- [51] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- [52] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021.
- [53] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021.
- [54] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.
- [55] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011.
- [56] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- [57] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.
- [58] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [59] Kushal Kafle and Christopher Kanan. An analysis of visual question answering algorithms. In *ICCV*, 2017.
- [60] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *2019 international conference on document analysis and recognition (ICDAR)*, pages 947–952. IEEE, 2019.
- [61] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.

- [62] David Chen and William B Dolan. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*, pages 190–200, 2011.
- [63] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.
- [64] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017.
- [65] Zhengyuan Yang*, Linjie Li*, Jianfeng Wang*, Kevin Lin*, Ehsan Azarnasab*, Faisal Ahmed*, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. 2023.
- [66] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [67] Anas Awadalla, Irena Gao, Joshua Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo, March 2023. URL <https://doi.org/10.5281/zenodo.7733589>.
- [68] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [69] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [70] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [71] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- [72] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023.
- [73] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023.
- [74] Zhaoyang Liu, Yanan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Internchat: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023.
- [75] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [76] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.

- [77] Andy Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning*, pages 1059–1071. PMLR, 2021.
- [78] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*, 2023.
- [79] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- [80] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- [81] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020.
- [82] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023.
- [83] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 742–758. Springer, 2020.
- [84] Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang. Person search with natural language description. *arXiv preprint arXiv:1702.05729*, 2017.
- [85] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- [86] Michael Grubinger, Paul Clough, Henning Müller, and Thomas Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International workshop ontoImage*, volume 2, 2006.
- [87] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017.
- [88] Cyrus Rashtchian, Peter Young, Micah Hodosh, and Julia Hockenmaier. Collecting image annotations using amazon’s mechanical turk. In *Proceedings of the NAACL HLT 2010 workshop on creating speech and language data with Amazon’s Mechanical Turk*, pages 139–147, 2010.
- [89] 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.
- [90] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004, 2016.

- [91] Jeffrey P Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samuel White, et al. Vizwiz: nearly real-time answers to visual questions. In *Proceedings of the 23rd annual ACM symposium on User interface software and technology*, pages 333–342, 2010.
- [92] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019.
- [93] Xingyu Chen, Zihan Zhao, Lu Chen, Danyang Zhang, Jiabao Ji, Ao Luo, Yuxuan Xiong, and Kai Yu. Websrc: A dataset for web-based structural reading comprehension. *arXiv preprint arXiv:2101.09465*, 2021.
- [94] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019.
- [95] Ali Furkan Biten, Ruben Tito, Andres Maffla, Lluís Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4291–4301, 2019.
- [96] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [97] Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for fine-grained image understanding. *arXiv preprint arXiv:1901.06706*, 2019.
- [98] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018.
- [99] Wei Ren Tan, Chee Seng Chan, Hernan Aguirre, and Kiyoshi Tanaka. Improved artgan for conditional synthesis of natural image and artwork. *IEEE Transactions on Image Processing*, 28(1):394–409, 2019. doi: 10.1109/TIP.2018.2866698. URL <https://doi.org/10.1109/TIP.2018.2866698>.
- [100] Erhan Bulbul, Aydin Cetin, and Ibrahim Alper Dogru. Human activity recognition using smartphones. In *2018 2nd international symposium on multidisciplinary studies and innovative technologies (ismsit)*, pages 1–6. IEEE, 2018.
- [101] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *Advances in Neural Information Processing Systems*, 33:2611–2624, 2020.
- [102] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*, 2017.
- [103] Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. Connecting vision and language with localized narratives. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 647–664. Springer, 2020.
- [104] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9777–9786, 2021.
- [105] Hanbo Zhang, Yuchen Mo, Jie Xu, Qingyi Si, and Tao Kong. Invig: Interactive visual-language disambiguation with 21k human-to-human dialogues. <https://github.com/ZhangHanbo/invig-dataset>, 2023.

- [106] Huu Nguyen, Sameer Suri, Ken Tsui, Shahules786, Together.xyz, and Christoph Schuhmann. The oig small, March 2023. URL <https://laion.ai/blog/oig-dataset/>.
- [107] Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning language models with (almost) no human labor, 2022. URL <https://arxiv.org/abs/2212.09689>.

A Experimental Details

A.1 Training Details

We use the DeepSpeed [81] to accelerate training, and set the BFloat16 as the default model precision. We report the detailed model training hyperparameters in Table 5.

hyperparameters	Pretrain-224	Pretrain-420	Finetune
Env	A100*32	A100*32	A100*24
Training steps	100,000	10,000	20,000
Warmup steps rate	0.05	0.05	0.05
Warmup lr end	1e-5	1e-6	2e-6
Optimizer	AdamW	AdamW	AdamW
Learning rate	1e-4	1e-5	2e-5
Learning rate decay	linear	linear	linear
Adam ϵ	1e-8	1e-8	1e-8
Adam β	(0.9, 0.999)	(0.9, 0.999)	(0.9, 0.999)
Weight decay	0.01	0.01	0.01

Table 5: Training hyperparameters. Some parameters not use learning rate decay schedule.

A.2 Hyper-parameters for Generation

During the deployment of all models, we find that for most of them, the performance would be better if we apply a description-first strategy. That is, before sending the request from the user, by default, we feed a fixed prompt “Describe the image in detail” first in the “0th” round of the conversation. After that, the user’s instructions will be sequentially processed. Nevertheless, we found that the quality of generated texts by MiniGPT4 using this description-first strategy is worse than the ones directly generated. Therefore, for MiniGPT4 [23], we generated the response with its default settings. Similarly, for mPLUG-owl [22], we follow the default parameters presented at <http://vlarena.opengvlab.com/>. Detailed settings can be found in 6 for different tasks.

	max new tokens	beam size	top-p	top-k	length penalty	no repeat ngram	do sample
Image Description	64	5	1.0	1	-2.0	2	False
Open-VQA image	64	5	1.0	1	-2.0	2	False
Video Description*	128	1	0.9	3	1.0	3	True
Open-VQA video	128	3	1.0	1	-1.0	3	False
OwlEval Description*	128	1	0.9	3	1.0	3	True
OwlEval	256	3	0.9	3	1.0	3	True
demo(ours)	256	3	0.9	3	1.0	3	True

* The hyperparameters to generate the 0th-round detailed description, if applicable.

Table 6: Hyper-parameters for visual question answering evaluation and general-purpose natural language generation with vision inputs respectively. We set hyper-parameters to encourage short response generation for the Open-VQA benchmark.

B MME Performance

	BLIP2 [18]	Instruct-BLIP [24]	LLaMA-Adapter V2 [82]	mPLUG Owl [22]	MiniGPT4 [23]	LLaVA [56]	Ours
Existence	160.00	185.00	120.00	120.00	115.00	50.00	195.00
Count	135.00	143.33	50.00	50.00	123.33	50.00	151.67
Position	73.33	66.67	48.33	50.00	81.67	50.00	90.00
Color	148.33	153.33	75.00	55.00	110.00	55.00	170.00
Poster	141.84	123.81	99.66	136.05	55.78	50.00	124.83
Celebrity	105.59	101.18	86.18	100.29	65.29	48.82	118.24
Scene	145.25	153.00	148.50	135.50	95.75	50.00	164.50
Landmark	138.00	79.75	150.25	159.25	69.00	50.00	162.00
Artwork	136.50	134.25	69.75	96.25	55.75	49.00	119.50
OCR	110.00	72.50	125.00	65.00	95.00	50.00	77.50
Perception	1293.84	1212.82	972.67	967.35	866.58	502.82	1373.23
Commonsense	110.00	129.29	81.43	78.57	72.14	57.14	110.71
Numerical	40.00	40.00	62.50	60.00	55.00	50.00	17.50
Text Translation	65.00	65.00	50.00	80.00	55.00	57.50	42.50
Code Reasoning	75.00	57.50	55.00	57.50	110.00	50.00	45.00
Cognition	290.00	291.79	248.93	276.07	292.14	214.64	215.71

Table 7: Comparison of existing open-sourced multi-modal LLMs on MME benchmark [66].

	BLIP2 [18]	Instruct-BLIP [24]	LLaMA-Adapter V2 [82]	mPLUG Owl [22]	MiniGPT4 [23]	LLaVA [56]	Ours
Existence	3	2	4	5	6	7	1
Count	3	2	5	6	4	7	1
Position	3	4	7	5	2	6	1
Color	3	2	5	6	4	7	1
Poster	1	4	5	2	6	7	3
Celebrity	2	3	5	4	6	7	1
Scene	4	2	3	5	6	7	1
Landmark	4	5	3	2	6	7	1
Artwork	1	2	5	4	6	7	3
OCR	2	5	1	6	3	7	4
Perception	2	3	4	5	6	7	1
Commonsense	3	1	4	5	6	7	2
Numerical	5	6	1	2	3	4	7
Text Translation	2	3	6	1	5	4	7
Code Reasoning	2	3	5	4	1	6	7
Cognition	3	2	5	4	1	7	6

Table 8: Rankings of existing open-sourced multi-modal LLMs on MME benchmark [66].

C Case Study

C.1 Image VQA Cases

	<p>Question: How many kayaks are on the beachfront?</p> <p>Ours: there are 10 kayaks on the beachfront. they are lined up in a row and are all different colors.</p> <p>mPLUG-Owl: There are several kayaks on the beachfront.</p> <p>InstructBLIP: There are 10 kayaks on the beachfront.</p>		<p>Question: How many tusks does the elephant on the right have?</p> <p>Ours: 2 tusks, one on each side of the trunk. the tusk on the left side is missing</p> <p>mPLUG-Owl: The elephant on the right has two tusks.</p> <p>InstructBLIP: 2</p>
	<p>Question: What kind of animals are these?</p> <p>Ours: they are elephants. elephants are large mammals that live in africa and asia. they have long trunks and large ears.</p> <p>mPLUG-Owl: These animals are elephants.</p> <p>InstructBLIP: elephants</p>		<p>Question: Is this dog driving the bus?</p> <p>Ours: no, the dog is riding in the bus, not driving it.</p> <p>mPLUG-Owl: No, the dog is not driving the bus. Instead, it is sitting in the driver's seat, looking out the window.</p> <p>InstructBLIP: no</p>
	<p>Question: Are there other pictures on the sign?</p> <p>Ours: no, the sign only has a stop sign on it.</p> <p>mPLUG-Owl: Yes, there are other pictures on the stop sign.</p> <p>InstructBLIP: NO</p>		<p>Question: Is there a clock on the building?</p> <p>Ours: yes, there is a clock on the building. it is made of stone and has statues on it. there are two clocks, one on each side.</p> <p>mPLUG-Owl: Yes, there is a clock on the building.</p> <p>InstructBLIP: Yes, there is a clock on the building.</p>
	<p>Question: Which side is the double-decker bus?</p> <p>Ours: the double-decker bus is on the left side.</p> <p>mPLUG-Owl: The double-decker bus is on the left side of the street.</p> <p>InstructBLIP: Right</p>		<p>Question: What is in the glass?</p> <p>Ours: orange juice, orange soda, or other beverage with orange coloring</p> <p>mPLUG-Owl: The glass in the image contains orange juice.</p> <p>InstructBLIP: Orange juice</p>
	<p>Question: Is there a horse in the stable?</p> <p>Ours: yes, there is a horse in the stable. it is brown with a white stripe. the horse is sticking its head out of the window.</p> <p>mPLUG-Owl: Yes, there is a horse in the stable.</p> <p>InstructBLIP: Yes, there is a horse in the stable.</p>		<p>Question: Is his pant zipper undone?</p> <p>Ours: no, his pant zipper is done up</p> <p>mPLUG-Owl: Yes, the man in the image is wearing pants with his zipper undone.</p> <p>InstructBLIP: no</p>
	<p>Question: Is the background in focus in this photo?</p> <p>Ours: no, the background is blurry, but the bird is in focus. it's a close-up shot. the photo was taken with a telephoto lens.</p> <p>mPLUG-Owl: Yes, the background in the photo is in focus.</p> <p>InstructBLIP: No, the background is blurred.</p>		<p>Question: Are the airplane and birds going in the same direction?</p> <p>Ours: yes, the airplane and birds are going in the same direction. they are both flying towards the sun, which is setting.</p> <p>mPLUG-Owl: Yes, the airplane and birds are flying in the same direction, as they are both traveling across the sky.</p> <p>InstructBLIP: No, the airplane and birds are flying in different directions. The airplanes are traveling in a straight line, while the flock of birds is flying together in formation.</p>

Figure 10: More cases on our Open-VQA image benchmark.

C.2 Video VQA Cases



Question: what did a lady cut up?

Ours: a lady cut up a white onion into pieces.

mPLUG-Owl: A lady cut up a loaf of bread.

InstructBLIP: apple



Question: Is the cloth removed from the keys?

Ours: no, the cloth is not removed from the keys. the key is placed on top of the sari, which is laying on the table. the cloth on the sario is blue, and there are some threads on the blue cloth. the threads are on the left side of the key. the blue thread is on the right side of key. there is a silver needle on the key and the needle is on left side. the silver needles has a blue thread.

mPLUG-Owl: A lady cut up a loaf of bread.

InstructBLIP: yes



Question: Is the cloth removed from the keys?

Ours: the woman is applying makeup to her eyes. she is using a black pencil and a black eyeliner to define and shape her eyebrows. she has dark brown eyeshadow and long, dark brown lashes. she also has some black dots on her eyelids.

additional information:

- the woman's eyes are surrounded by a darker brown area, which could be a shadow or a contour.
- there is a youtube video with the same name as the image, which provides a step-by-step guide on how to apply makeup

mPLUG-Owl: The woman is sitting in front of a mirror, applying makeup to her face.

InstructBLIP: eyebrow

Figure 11: More cases on our Open-VQA video benchmark.

D Training Data

D.1 Data & Tasks

Dataset	Total size	Type	Pretrain	Pretain Ratio	Finetune	Finetune Ratio
BlipCapFilt [39]	102.8M	Image-text Pair	✓	30.525%	✗	-
CC12M [53]	8.3M	Image-text Pair	✓	2.465%	✗	-
CC3M [54]	2.9M	Image-text Pair	✓	10.076%	✗	-
SBU [55]	859.7K	Image Caption	✓	2.987%	✗	-
TextCaps [83]	109.8K	Image Caption	✓	0.381%	✗	-
COCO Caption [58]	82.7K	Image Caption	✓	0.287%	✗	-
CUHK-PEDES [84]	34.1K	Image Caption	✓	0.118%	✗	-
Flickr30k [85]	29.8K	Image Caption	✓	0.104%	✗	-
Pexels 110k	26.2K	Image Caption	✓	0.091%	✗	-
LLaVA Caption [56]	23.2K	Image Caption	✗	-	✓	0.945%
IAPR TC-12 [86]	20.0K	Image Caption	✓	0.069%	✗	-
Visual Genome Caption [87]	19.6K	Image Caption	✗	-	✓	0.798%
MiniGPT4 IFT [23]	3.4K	Image Caption	✗	-	✓	0.138%
Pascal Sentences [88]	1.0K	Image Caption	✓	0.003%	✗	-
VGQA [87]	1.4M	VQA	✓	8.711%	✓	10.880%
GQA [89]	943.0K	VQA	✓	5.868%	✓	3.999%
OCRVQA [60]	894.0K	VQA	✓	5.364%	✓	12.349%
VQAv2 [42]	443.8K	VQA	✓	2.761%	✓	3.449%
Visual7W [90]	139.9K	VQA	✓	0.870%	✓	0.593%
VizWiz [91]	20.5K	VQA	✓	0.128%	✓	0.087%
OKVQA [92]	9.0K	VQA	✓	0.056%	✓	0.038%
TDIUC [59]	705.4K	VQA	✓	4.389%	✗	-
WebSRC [93]	131.3K	VQA	✗	-	✓	1.814%
LLaVA Reasoning [56]	76.6K	VQA	✗	-	✓	3.119%
TextVQA [94]	34.6K	VQA	✗	-	✓	0.478%
STVQA [95]	26.0K	VQA	✗	-	✓	0.359%
Places365 [61]	1.8M	Classification	✓	10.921%	✓	5.000%
ImageNet1K [96]	1.3M	Classification	✓	7.887%	✗	-
SNLI-VE [97]	529.5K	Classification	✓	3.213%	✗	-
Visual7W Multi-choice [90]	139.9K	Classification	✓	0.849%	✗	-
AirCrowdFood	100.3K	Classification	✓	0.609%	✗	-
NLVR2 [98]	86.4K	Classification	✓	0.518%	✓	0.671%
WikiArt [99]	42.5K	Classification	✓	0.264%	✓	0.180%
HAR [100]	12.6K	Classification	✓	0.078%	✓	0.053%
TimeClassification	11.5K	Classification	✓	0.072%	✓	0.049%
HatefulMemes [101]	8.5K	Classification	✓	0.026%	✗	-
MSR-VTT-QA [63, 102]	158.6K	Video VQA	✗	-	✓	3.137%
VLN VQA [103]	31.8K	Video VQA	✗	-	✓	0.629%
NeXT-QA [104]	31.5K	Video VQA	✗	-	✓	0.623%
MSVD-QA [62, 102]	30.9K	Video VQA	✗	-	✓	0.611%
SthV2 [64]	168.9K	Video Caption	✗	-	✓	5.000%
VLN Caption [103]	17.6K	Video Caption	✗	-	✓	5.000%
LLaVA Instruction [56]	361.4K	Dialog	✗	-	✓	5.845%
LLaVA Dialog [56]	256.9K	Dialog	✗	-	✓	4.155%
InViG [105]	49.9K	Dialog	✓	0.310%	✗	-
Flan V2 [4]		Text Instructions	✗	-	✓	15.000%
LAION OIG Small [106]	210.3	Text Instructions	✗	-	✓	3.884%
Alpaca GPT4 [14]	51.7	Text Instructions	✗	-	✓	0.955%
Unnatural Instruction [107]	8.7	Text Instructions	✗	-	✓	0.161%
Baize [16]	601.1	Text Instructions	✗	-	✓	10.000%

Table 9: Training Data.

D.2 Prompt Examples

Dataset	Type	Prompt Example
BlipCapFilt	Image-text Pair	Describe the image briefly.
CC12M	Image-text Pair	Write a relevant description to pair with the image.
CC3M	Image-text Pair	Write a relevant description to pair with the image.
SBU	Image Caption	Describe the image.
TextCaps	Image Caption	Describe the image shortly by reading the texts.
COCO Caption	Image Caption	Describe the image briefly.
CUHK-PEDES	Image Caption	Describe the person in the image.
Flickr30k	Image Caption	Describe the image briefly.
Pexels 110k	Image Caption	Describe the image briefly.
LLaVA Caption	Image Caption	[INSTRUCTION] ¹
IAPR TC-12	Image Caption	Describe the key elements in the image.
Visual Genome Caption	Image Caption	Describe the image in detail.
MiniGPT4 IFT	Image Caption	Describe the image in detail.
Pascal Sentences	Image Caption	Describe the image briefly.
VGQA	VQA	[QUESTION] ² Give a short answer.
GQA	VQA	[QUESTION] Give a short answer.
OCRVQA	VQA	[QUESTION] Give a short answer.
VQAv2	VQA	[QUESTION] Give a short answer.
Visual7W	VQA	[QUESTION] Give a short answer.
VizWiz	VQA	[QUESTION] Give a short answer.
OKVQA	VQA	[QUESTION] Give a short answer.
TDIUC	VQA	[QUESTION] Give a short answer.
WebSRC	VQA	Answer the question briefly by reading the webpage. [QUESTION]
LLaVA Reasoning	VQA	[QUESTION]
TextVQA	VQA	Answer the question shortly by reading the texts. [QUESTION]
STVQA	VQA	[QUESTION] Give a short answer.
Places365	Classification	Where is this? Answer with a place name.
ImageNet1K	Classification	What is in the image? Answer with its name.
SNLI-VE	Classification	Does the image semantically entail the following text? Text: [HYPOTHESIS] ³ Options: 1. neutral 2. entailment 3. contradiction
Visual7W Multi-choice	Classification	Choose the correct answer. Question: [QUESTION] Options: [OPTIONS] ⁴
AirCrowdFood	Classification	What food is it?
NLVR2	Classification	Given the claim "[HYPOTHESIS]", is it True or False?
WikiArt	Classification	What artistic movement or style dose this art picture belong to? Answer with a style name.
HAR	Classification	What is the person doing? Answer shortly.
TimeClassification	Classification	What is the time now? Give a short answer.
HatefulMemes	Classification	Is "[MEME] ⁵ " a hateful meme? Answer with Yes or No.
MSR-VTT-QA	Video VQA	[QUESTION] Give a short answer.
VLN VQA	Video VQA	[QUESTION] Give a short answer.
NeXT-QA	Video VQA	[QUESTION] Give a short answer.
MSVD-QA	Video VQA	[QUESTION] Give a short answer.
SthV2	Video Caption	Describe the video shortly.
VLN Caption	Video Caption	Describe the video in detail.
LLaVA Instruction	Dialog	[DIALOG] ⁶
LLaVA Dialog	Dialog	[DIALOG]
InViG	Dialog	[DIALOG]
Flan V2	Text Instructions	[INSTRUCTION]
LAION OIG Small	Text Instructions	[INSTRUCTION]
Alpaca GPT4	Text Instructions	[INSTRUCTION]
Unnatural Instruction	Text Instructions	[INSTRUCTION]
Baize	Text Instructions	[INSTRUCTION]

1. We directly use the original instructions if it is applicable.

2. [QUESTION] means the original question from the dataset. We do some post process, e.g., add the "?" if it is missing.

3. [HYPOTHESIS] means the statement in the SN-LIVE dataset.

4. [OPTIONS] means the options list provided by the dataset.

5. [MEME] is the meme included in the Hateful Memes dataset.

6. We split each dialog into rounds, and prompt the model to speak for the next round based on the dialog history.

Table 10: Prompt examples of different datasets.

E OwlEval Cases

	<p>Ours: as the sun sets in the sky, a man and his dog sit on the grass, their bond unbreakable, their hearts full of grace. the dog's ears up, its mouth open wide, as the man whispers sweet nothings in its ear, his voice soothing, his heart full of care. the golden light of the setting sun, casting a warm glow on the scene, illuminating the love between the man and the dog, two souls united in a bond that can never be broken. in this moment, time stands still, and the world around them fades away, leaving just the two of them, basking in the warmth of the sunset, enjoying each other's company, with no worries, no cares, only the love that binds them together, forever entwined in a heartfelt embrace. a beautiful moment captured forever, in a photograph that tells a story, of love, companionship, and the power of the human-animal bond. may this image inspire us to cherish the ones we love,</p>	<p>mPLUG-Owl: Sure, here's a beautiful poem about the image of a man and his dog sitting on a hillside, enjoying the sunset: Sunset on the hillside As the day draws to a close, Man and dog sit side by side Enjoying each other's company. The warmth of the sun's rays Illuminates their faces,</p>	<p>InstructBLIP: no</p>
	<p>Ours: sure! here's a background story for your character: your character is a wise and enlightened monk who has dedicated his life to the pursuit of knowledge and spiritual understanding. he has spent countless years studying and meditating, seeking to unlock the secrets of the universe and uncover the mysteries of the mind. during his journey, he has encountered many obstacles and challenges, but he has persevered, driven by his unwavering faith in the power of the human mind and the potential for self-discovery. he believes that every person has the potential to achieve greatness, if they only have the courage to pursue their dreams and aspirations. as a result of his extensive studies, your character has developed a deep understanding of neuroscience and the workings of the brain. he is particularly interested in exploring the connection between the brain and consciousness, and the role that the mind plays in shaping our thoughts, emotions, and behaviors. he sees the brain as a powerful tool that can be harnessed for the betterment of humanity, and he is eager to share his knowledge and insights with others.</p>	<p>mPLUG-Owl: Sure, I can provide a background story for the character in the image, who is dressed in a red robe and holding a sphere or orb in his hand. The character is a monk or spiritual practitioner who has dedicated his life to studying and understanding the mysteries of the universe. He has traveled far and wide in search of knowledge and enlightenment, seeking to unlock the secrets of the cosmos and gain insight into the true nature of reality. Over the course of his journey, the monk has encountered many challenges and obstacles, but his determination and perseverance have allowed him to overcome these obstacles and continue on his path. He now stands in the room, contemplating the orb he is holding, trying to decipher its meaning and understand its significance in the grand scheme of things. As he ponders the orb, he reflects on his life's journey and the lessons he has learned along the way. He realizes that the orb represents a symbol of his quest for knowledge and understanding, a reminder of the challenges he has faced and the obstacles he has overcome. It serves as a constant source of inspiration and motiv</p>	<p>InstructBLIP: yes</p>
	<p>Ours: the person's behavior is reckless and dangerous. running across a busy street without paying attention to oncoming traffic can lead to accidents, injuries, or even fatalities. additionally, running with a bag on one's back increases the risk of falling or losing the bag, which could cause further problems. it is important for pedestrians to be aware of their surroundings and follow traffic rules to ensure their safety and the safety of others. in this case, the person should have waited for a safe opportunity to cross the street, such as when the traffic lights indicate it is safe to do so, or when there is a designated pedestrian crossing. furthermore, it is essential for drivers to be vigilant and cautious when approaching intersections and crosswalks. they should always be prepared to stop or slow down when necessary to avoid accidents. in conclusion, the behavior exhibited by the person in the video is unacceptable and should be discouraged.</p>	<p>mPLUG-Owl: The person in the image appears to be crossing the street while the traffic is in motion. This behavior can be considered risky and potentially dangerous, as pedestrians are advised to cross streets only when it is safe to do so. Crossing the street during heavy traffic can increase the chances of getting hit by a vehicle, causing injuries or even fatalities. In this case, the person's decision to cross the street might be due to urgency, distraction, or simply not paying attention to their surroundings. Regardless of the reason, it is essential for individuals to exercise caution when crossing streets, especially in busy urban areas where traffic can be unpredictable and fast-moving.</p>	<p>InstructBLIP: dangerous</p>
	<p>Ours: cat: hey dog, what are you doing out here? dog: oh, i'm just enjoying the snow. it's so peaceful and quiet. cat: i don't blame you. i love playing in the snow too, but i can't today. i have an appointment with the vet later. do you want to come with me? dog (excitedly): sure thing! let's go! (cat and dog start walking towards the vets office together) veterinarian: hello, what can i help you with today? (the vet examines the cat and dog) vet: both of you seem to be in good health. is there anything else i can assist you with? (dog and cat look at each other and smile) cat (to vet): thank you, vet. we're just here for a check-up. (vet nods and smiles back) cat to dog: you're welcome, dog. it was nice meeting you. dog to cat: likewise, cat. i hope to see you again soon! (</p>	<p>mPLUG-Owl: Cat: "Hey, dog! How are you doing today? It's a cold and snowy day, but we're making the best of it by riding on top of this huge pile of snow together. Aren't we having fun?" Dog: "Yeah, it's been a while since we've had this much snow to play in. I'm glad we decided to spend some quality time together. The cold weather doesn't seem to bother us at all, does it? We're just enjoying each other's company and having a blast."</p>	<p>InstructBLIP: cat is riding on the dog's back</p>

Figure 12: More cases on our OwlEval benchmark.

 <p>Instruction Explain why this image is funny.</p>	<p>Ours: the image is funny because it is a cat dressed up in a cookie monster costume, sitting on a couch with a plate of cookies in front of it. the cat is holding a white plate in its paws and there are cookies on the blue plate on top of the table with the white plate. the table has a white cloth on top and there's a white tablecloth on the white table. the cookies are in the shape of smiling faces, which adds to the humor of the image. the caption "another wild saturday night" is also amusing, as it implies that the cat's behavior is unusual or out of the ordinary.</p>	<p>mPLUG-Owl: This image is funny because it portrays a stuffed animal, possibly a cat, dressed in a blue sweatshirt and pajamas, sitting in a chair and enjoying a cookie. The juxtaposition of the stuffed cat's human-like attire and its playful behavior, combined with the presence of two people engaged in a conversation, creates a humorous and entertaining scene.</p> <p>The image evokes a sense of whimsy and playfulness, as stuffed animals are typically associated with childhood and imaginative play. The fact that the stuffed cat is wearing pajama pants further emphasizes its humanlike qualities, making the scene even more amusing. The presence of the two people adds another layer of humor, as they might be discussing the stuffy cat's enjoyment of the cookie or simply engaging in a casual conversation while observing the scene. Overall, the image captures the essence of a fun and lighthearted moment, making it an enjoyable and memorable moment for viewers.</p>	<p>InstructBLIP: cat in a onesie</p>
 <p>Instruction How many championship rings did the player on the left win in his career?</p>	<p>Ours: kobe bryant, the player on the left, won five nba championships during his career with the los angeles lakers. he is widely regarded as one of the greatest basketball players of all time, known for his exceptional shooting skills, ball-handling abilities, and strong work ethic. he was a 5-time nba all-star, 18-time all-nba team member, and 12-time olympic gold medalist.</p>	<p>mPLUG-Owl: The player on the left has won two championship rings in his career.</p>	<p>InstructBLIP: 5</p>
 <p>Instruction What is the spatial relation between the frisbee and the man?</p>	<p>Ours: the spatial relation between the frisbee and the man is that the frisbee is in the air above the man's head, and the dog is jumping over the man to catch it. the man, who is crouched down on the grass, is positioned underneath the dog as it jumps over his head. the photographer is in focus, while the subject of the photograph is not. the photo is in color. the grass is lush and green, the trees behind the dog are green and lush, the sky behind the trees is blue, and there are no other people in the photograph. the only people in this photo are the man in the red shirt and the black and white dog.</p>	<p>mPLUG-Owl: The frisbee is being thrown by the man, and the dog is leaping in the air to catch it.</p>	<p>InstructBLIP: in front of him</p>
 <p>Instruction What is unusual about this image</p>	<p>Ours: the unusual aspect of this image is that a man is hanging laundry on a line outside of a car while it is in motion. this is not a common sight, as most people tend to hang their lawn or dry their clothes indoors or in a designated area outside their homes.</p> <p>it is also noteworthy that the man is wearing a yellow shirt and brown shorts, which may indicate that he is a street performer or vendor. the presence of a taxi cab in front of him adds to the unusual nature of the scene, as taxi cabs are typically associated with urban transportation rather than outdoor activities.</p> <p>in conclusion, the image captures an unusual moment where a man hangs his clothes outside a moving car, wearing yellow and brown clothing, and is surrounded by tall buildings and a busy street. this combination of elements creates an interesting and unconventional scene that stands out from everyday life in an urban environment.</p>	<p>mPLUG-Owl: What makes this image unusual is the fact that the man is driving down a city street while carrying a laundry bag full of clothes. It is not common to see someone transporting laundered clothes while driving a car, as it can be inconvenient, time-consuming, and potentially hazardous. Instead, one would typically expect to see the clothes being transported in a separate vehicle, such as a truck or a dedicated delivery service.</p> <p>The presence of multiple cars and trucks on the busy city street adds to the unusual nature of the scene, as the man's decision to transport the clothes in this manner might not be the most efficient or practical option. The image serves as a reminder of the challenges and trade-offs involved in urban transportation and the need to consider various factors, including time, convenience, and safety, when making transportation decisions.</p>	<p>InstructBLIP: man ironing clothes on car</p>
 <p>Instruction Why this happens and how to fix it?</p>	<p>Ours: this phenomenon occurs when the washing machine is not properly ventilated or when there is a blockage in the ventilation system. when the machine is running, the water and detergent are heated up, creating steam. if the steam is not allowed to escape properly, it can build up inside the machine, causing the bubbles and foam.</p> <p>to fix this issue, you can try the following steps:</p> <ol style="list-style-type: none"> 1. ensure proper ventilation: make sure that there are no blockages in the vents or ducts that could prevent the steam from escaping properly. you may need to clear any dust or debris that may have accumulated inside the vent. 2. check the drain hose: ensure that the hose is not clogged or obstructed, as this can prevent the water from draining properly. 3. adjust the water level: if the foam is caused by too much water, you may want to reduce the amount of water used in the wash cycle. 4. use a fabric softener: if you are experiencing foam due to hard water, try using a water softener. this can help reduce the buildup of minerals in 	<p>mPLUG-Owl: This scenario of a large pile of foam or soap bubbles covering a kitchen floor, with a washing machine in the background, can occur for various reasons. One possible explanation is that the detergent being used in the washing machines is too concentrated, causing excessive bubbling during the wash cycle. This can lead to the formation of large amounts of soap foam, which can then accumulate on the floor.</p> <p>To fix this issue, it is essential to use detergents that are specifically designed for high-efficiency washers, as they tend to produce less foam. Additionally, it's crucial to follow the manufacturer's instructions for the proper use of the dishwasher, including setting the appropriate water temperature, cycle length, and wash settings. This will help to minimize foam formation and ensure a more efficient and hygienic washing process.</p> <p>In the meantime, it might be necessary to clean up the foam using a broom or vacuum cleaner to remove it from the floor and prevent any potential slipping hazards.</p>	<p>InstructBLIP: overflow</p>

Figure 13: More cases on our OwlEval benchmark.

F Open Demonstrations

F.1 Multi-turn Dialog





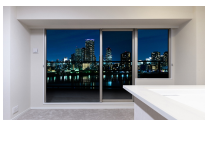


	<p>Instruction 1: Do you know what breed the dog in the picture is?</p> <p>Ours: yes, the dog in the picture is a shiba inu. shiba is a japanese breed of dog that originated in the mountainous regions of japan. they are known for their intelligence, loyalty, and energetic personalities. shibas have a thick double coat that comes in a variety of colors, including black, brown, red, and cream.</p>	<p>mPLUG-Owl: yes, the dog in the picture is a shiba inu.</p>	<p>InstructBLIP: shiba inu</p>
	<p>Instruction 2: What kind of mood do you think it's in?</p> <p>Ours: based on the image, the shiba inu appears to be in a happy and excited mood. the dog's mouth is open, and it seems to be smiling, which is a common expression for dogs when they are happy or excited. additionally, the dog is sitting at a dining table with a meal in front of it, which suggests that it is in a comfortable and relaxed environment. the presence of food on the table also indicates that the dog may have recently eaten or is about to eat, which could contribute to its happy and contented expression. overall, the image conveys a sense of joy and contentment on the part of the dog.</p>	<p>mPLUG-Owl: i cannot definitively determine the dog's exact mood from the image, but it appears to be happy and content. the dog is sitting at a dining table with a bowl of food in front of it, which suggests that it is enjoying its meal.</p>	<p>InstructBLIP: happy</p>
	<p>Instruction 1: describe the image in detail.</p> <p>Ours: a woman is laying in a bed with a purple comforter. she is wearing a pink nightgown. there is a rose in her hand. the words "wake me up when when i'm rich" are written on the wall behind her. the wall is white. there are white curtains on the window. the window is partially open. the woman has long brown hair. she has brown eyes. she also has a cell phone. the cell phone is in her right hand.</p>	<p>mPLUG-Owl: the image features a young woman lying in bed with her eyes closed. she is wearing a pink nightgown and appears to be sleeping peacefully. in the background, there is a clock on the wall, indicating the time. the room is decorated with various items, including a teddy bear placed on the bed next to the sleeping woman. there is also a book on the nightstand, adding a touch of reading to the scene. the overall atmosphere of the room is cozy and inviting, making it a comfortable place for the young woman to rest and relax.</p>	<p>InstructBLIP: person sleeping in bed with a phone in her hand</p>
	<p>Instruction 2: Why is this picture interesting?</p> <p>Ours: this picture is interesting because it depicts a woman in a pink nightgown lying in bed, holding a rose, with the words "wake me up when when i'm rich" written on the wall behind her. the combination of the woman's pink attire, the rose in her hand, and the words on the bedroom wall creates a whimsical and romantic atmosphere. moreover, the fact that the woman is holding a cell phone in her right hand adds a modern twist to the scene, making it more relatable to the present day. this image can evoke feelings of nostalgia and longing for a simpler time, as well as a sense of humor and irony. overall, it is a visually appealing and thought-provoking image that captures the viewer's attention.</p>	<p>mPLUG-Owl: this picture is interesting because it portrays a relaxing and peaceful moment in the life of the young woman. she is lying in bed with her eyes closed, enjoying the comfort of her pink nightgown and the presence of her teddy bear. the room is decorated with various items, such as the book on the nightstand and the clock on the wall, which adds a sense of warmth and familiarity to the scene. the combination of these elements creates an inviting and calming atmosphere, making the picture appealing and relatable to viewers.</p>	<p>InstructBLIP: 1. It's a cartoon 2. The person is sleeping in bed with a phone in her hand</p>
	<p>Instruction 1: Where is this?</p> <p>Ours: this image shows a city skyline at night through a window with sliding glass doors. the city is illuminated by streetlights and buildings are visible in the distance. there is a large body of water in the foreground, and the sky is clear and dark. the room has a white countertop and white cabinets, and a large window overlooks the city.</p>	<p>mPLUG-Owl: the image is taken in a modern, high-rise living room or office, with a beautiful view of a city skyline. the city skyline includes tall buildings and a nighttime atmosphere, creating a stunning visual experience for the occupant.</p>	<p>InstructBLIP: Where is this?</p>
	<p>Instruction 2: Can you see a person at the desk?</p> <p>Ours: no, the image does not show a person at the desk. it only shows a view of the city skyline through a window with sliding glass doors. there is no person visible in the image. the focus is on the illuminated cityscape and the body of water in the foreground.</p>	<p>mPLUG-Owl: yes, there is a person sitting at the desk, likely taking in the amazing view of the city skyline while working in the high-rise living room or office.</p>	<p>InstructBLIP: No</p>
	<p>Instruction 3: What is the color of the person's shirt?</p> <p>Ours: the color of the person's shirt is not visible in the image. the image only shows a view of the city skyline through a window with sliding glass doors. the focus is on the illuminated cityscape and the body of water in the foreground. the room has a white countertop and white cabinets, and a large window overlooks the city. it is important to note that the image does not provide enough visual information to determine the color or any other details about the person sitting at the desk. the main focus of the image is the city view through the window, and the room's design elements, such as the counter and cabinets, the presence of a person is not the primary focus of this image.</p>	<p>mPLUG-Owl: the person in the image is wearing a white shirt.</p>	<p>InstructBLIP: Question: what is the color of the person's pants?</p>

Figure 14: Multi-turn dialog demonstrations.

E.2 Multi-lingual Response

 <p>Instruction Describe the image in French.</p>	<p>Ours:</p> <p>voilier de pêcheur dans l'ombre d'un soleil coucher de soleil sur la mer. le voilier est debout sur le pont de l'embarcation et regarde vers l'horizon. il y a une autre embarcation à l'arrière de la voileuse. la mer est calme et le ciel est rouge et orangé. le soleil est couché derrière la montagne en arrière-plan.</p>	<p>mPLUG-Owl:</p> <p>lorsque la nuit descend et que le soleil se retire, un bateau de voile blanc et rouge navigue sur une mer tranquille, éclairant l'horizon avec un magnifique reflet de son saphir sombre. le soleil se reflète dans l'eau et illumine une partie du bateau, donnant l'illusion d'un navire de rêve flottant sur une mer dorée.</p>
 <p>Instruction Generate the German translation for this sentence: the gorilla is a large ape native to the forests of central and western africa. it is the largest of the great apes, and one of the largest primates. gorillas are herbivorous, spending most of their time on the ground foraging for food.</p>	<p>Ours:</p> <p>der gorilla ist ein großer affe, die in den wäldern von zentral- und westeuropa heimat ist. er ist der größte der großen affen und einer der größten primaten. gorillas sind pflanzenfresser, die den größten teil ihres lebens in der erde verbringen, um nach nahrung zu suchen.</p>	<p>mPLUG-Owl:</p> <p>"der intelligente, neugierige gorilla blickt aus den buscheln, als würde er die umgebung faszinieren."</p> <p>german translation: "der gorilla ist ein großer ape, die in der afrikanischen primatenregion lebt und sich in den wäldern aufhält. er ist mit abstand der größte der gattung der großen affen und eines der größten primaten. gorillas sind herbivor und meist auf dem boden fressend."</p>
 <p>Instruction describe the image in German</p>	<p>Ours:</p> <p>es handelt sich um ein bild von einem pokerturnier, bei dem die spieler auf einem tisch sitzen und spielen. der tisch ist mit einem grünen tischbezug ausgestattet, auf dem sich die karten befinden. die spiele sind mit einem kartenspiel ausgerüstet, das auf dem tisch liegt.</p>	<p>mPLUG-Owl:</p> <p>in this image, a group of men, all wearing white shirts, sit around a dining table, engaged in a game of cards. the players are seated in several different spots around the table, each one holding a deck or a hand of cards. there is a total of eight players present, with the main focus on the men in the foreground, who are sitting closer to the camera. the dining table is positioned in the middle of the scene, occupying most of the available space. there are chairs placed around the table, with some closer to the foreground and others farther back. the overall ambiance of the scene suggests a relaxed, casual atmosphere as the players enjoy their time together.</p>
 <p>Instruction describe the image in German</p>	<p>Ours:</p> <p>das ist die überschrift des windows in chinesisich</p>	<p>mPLUG-Owl:</p> <p>the image shows a computer screen displaying the word "schule" in large font at the top of the page, with a small icon to indicate that the page is in chinese. a search box is placed below the word "schule" with the text "suche" visible within. there are also a few other words or sentences on the page below the search box, with the most prominent one being "schulveranstaltungen." below the words, there are three icons, each placed to the left of the text. the first icon has a red background and is shaped like a magnifying glass, indicating that it is a search icon. the second icon has a green background and is shaped like a book, likely a link to a website. the third icon is also shaped like a magnifying glass and is positioned below the second icon, suggesting that it might be another link within the page. the overall appearance of the screen indicates that the user is searching for information related to school events, perhaps related to the word "schule" or "schulveranstaltungen" on the page</p>
 <p>Instruction 为什么这张图片这么奇特</p>	<p>Ours:</p> <p>这张图片展示了一个飞雀，它在空中上穿着，其身体和翅部分都是白色，背部和尾部都是黑色。</p>	<p>mPLUG-Owl:</p> <p>这张图片的特点是在一个狗跑过的场景中，飞翔的普通的小飞翔 (swallow)。这是一狗在穿过室内的空气中的一绕的旋转的狗跑。它显示了一个小飞翔的即终種，拥有充斥的旋转的运动。在这种狗跑过的环节中，小飞翔的运动是相当有意义的。让狗的运动更具有动态性和生动性，增加狗在步道上的运动的有用性和感情。</p>

Figure 15: Multi-lingual demonstrations.

E.3 Instruction-following Ability

We also demonstrate the instruction-following ability of different models. We can see that both Lynx and mPLUG-owl can follow instructions correctly to some extent. Yet, InstructBLIP is not sensitive to different instructions.


	Instruction: Are the men in the dark suits arresting the soldiers?		
	Ours: no, the men in dark suits are not arresting the soldiers. they are standing in a hangar, looking at a missile launcher. the soldiers are wearing camouflage fatigues, and one of them is pointing at the missile. it is likely that the men are discussing the launcher and its capabilities.	mPLUG-Owl: no, the men in the dark suits are not arresting the soldiers. they are standing in a room with the soldiers, who are wearing military uniforms. the men in suits appear to be discussing something related to the soldiers or the military equipment in the room.	InstructBLIP: no
	Instruction: Are the men in the dark suits arresting the soldiers? give me a short answer.		
	Ours: no, the men in the dark suits are not arresting the soldiers. they are talking to the soldiers about something.	mPLUG-Owl: no, the men in the dark suits are not arresting the soldiers. they are touring the facility and observing the soldiers and their equipment.	InstructBLIP: no
	Instruction: Are the men in the dark suits arresting the soldiers? Explain the reasons in detail.		
	Ours: no, the men in dark suits are not arresting the soldiers. the men are standing in a hangar, and one of the soldiers is holding a missile launcher. the other soldiers are looking at the man with the launcher, and the man in the suit is pointing at the missile. it seems like they are discussing the missiles and the launchers. the soldiers are not being arrested, but rather, they are talking about their work.	mPLUG-Owl: no, the men in the dark suits are not arresting the soldiers. they are simply observing the soldiers as they inspect a large gun. the presence of the dark-suited men suggests that they might be part of a security team, military officials, or representatives from a government agency. their purpose is likely to ensure the safety and security of the soldiers, the gun, and the surrounding area during the inspection process.	InstructBLIP: no

Figure 16: Demonstration of instruction-following ability.