# EMPOWERING VISION-LANGUAGE MODELS TO FOLLOW INTERLEAVED VISION-LANGUAGE INSTRUCTIONS

**Juncheng Li**[1,2*] **Kaihang Pan**[1*] **Zhiqi Ge**[1*] **Minghe Gao**[1*] **Hanwang Zhang**[3]
**Wei Ji**[2] **Wenqiao Zhang**[1] **Tat-Seng Chua**[2] **Siliang Tang**[1†] **Yueting Zhuang**[1†]

[1]Zhejiang University, [2]National University of Singapore, [3]Nanyang Technological University

## ABSTRACT

Multimodal Large Language Models (MLLMs) have recently sparked significant interest, which demonstrates emergent capabilities to serve as a general-purpose model for various vision-language tasks. However, existing methods mainly focus on limited types of instructions with a single image as visual context, which hinders the widespread availability of MLLMs. In this paper, we introduce the I4 benchmark to comprehensively evaluate the instruction following ability on complicated interleaved vision-language instructions, which involve intricate image-text sequential context, covering a diverse range of scenarios (*e.g.,* visually-rich webpages/textbooks, lecture slides, embodied dialogue). Systematic evaluation on our I4 benchmark reveals a common defect of existing methods: the Visual Prompt Generator (VPG) trained on image-captioning alignment objective tends to attend to common foreground information for captioning but struggles to extract specific information required by particular tasks. To address this issue, we propose a generic and lightweight controllable knowledge re-injection module, which utilizes the sophisticated reasoning ability of LLMs to control the VPG to conditionally extract instruction-specific visual information and re-inject it into the LLM. Further, we introduce an annotation-free cross-attention guided counterfactual image training strategy to methodically learn the proposed module by collaborating a cascade of foundation models. Enhanced by the proposed module and training strategy, we present **Cheetor**, a Transformer-based MLLM that can effectively handle a wide variety of interleaved vision-language instructions and achieves state-of-the-art zero-shot performance across all tasks of I4, without high-quality multi-modal instruction tuning data. Moreover, Cheetor also exhibits competitive performance compared with state-of-the-art instruction tuned models on concurrent MME benchmark. Our benchmark, code, and pre-trained models are available at `https://github.com/DCDmllm/Cheetah`.

## 1 INTRODUCTION

Large language models (LLMs) (OpenAI, 2023a;b; Touvron et al., 2023a;b; Chiang et al., 2023) have recently exhibited remarkable abilities on serving as a general-purpose model for wide-ranging tasks through instruction tuning (Wei et al., 2021; Wang et al., 2022) on collections of instructional style language tasks. By fine-tuning on various tasks in a unified instruction-response format, instruction tuning unlocks significant zero-shot generalizability of LLMs on novel task instructions. This success has inspired a new wave of research on extending text-only instruction-following models to multi-modal ones, with a longstanding aspiration in various real-world applications.

To achieve this goal, Flamingo (Alayrac et al., 2022) and BLIP-2 (Li et al., 2023b) empower LLMs with a frozen visual encoder to understand visual inputs. Follow-up works of LLaVA (Liu et al., 2023), MiniGPT-4 (Zhu et al., 2023), InstructBLIP (Dai et al., 2023) further improve the instruction-following ability of models by fine-tuning on multi-modal instruction following datasets. While these Multimodal Large Language Models (MLLMs) demonstrate inspiring capabilities, they mainly

---

*Equal Contribution. †Corresponding Authors.

Figure 1: Demonstrations and task taxonomy of the proposed I4 benchmark.

focus on vision-language instructions that only involve **a single image as visual context** and **limited instruction diversity**, thus restricting the availability of such instruction-following assistants. However, in real life, humans tend to convey their demands through a sequence of relevant images and texts. For example, people might require models to answer an open-domain question by referring to multiple sources of multi-modal information (*e.g.,* visually-rich webpages, textbooks, lecture slides). These diverse references and the question constitute **interleaved vision-language instructions**, where multiple images and texts are semantically interconnected.

To facilitate research in interleaved vision-language instruction following, we build **I4** (semantically **I**nterconnected, **I**nterleaved **I**mage-Text **I**nstruction-Following), an extensive large-scale benchmark of 31 tasks with diverse instructions in a unified instruction-response format, covering 20 diverse scenarios. As illustrated in Figure 1, I4 has three important properties: **(1) Interleaved vision-language context,** all the instructions contain sequences of inter-related images and texts, such as storyboards with scripts, and textbooks with diagrams. **(2) Diverse forms of complex instructions,** the instructions range from predicting dialogue for comics, to discovering differences between

2

surveillance images, and to conversational embodied tasks. **(3) Vast range of instruction-following scenarios,** the benchmark covers multiple practical scenarios, including cartoons, industrial images, driving recordings, recipes, *etc*.

Based on the proposed benchmark, we systematically evaluate recent MLLMs and find that they struggle to accomplish such complex multi-modal instructions. While existing MLLMs mainly focus on designing advanced methods to construct more diverse and high-quality instruction tuning data, we argue that **the Visual Prompt Generator (VPG) matters significantly in MLLMs for comprehending complex instructions**. To adapt LLMs to understand visual inputs, existing methods propose several VPGs (*e.g.,* linear projection, Resampler, Q-former) to extract relevant visual prompts from the dense image features encoded by the vision backbones (*e.g.,* ViT). They train the VPG on millions of image-caption pairs by requiring the frozen LLM to generate captions conditioned on the visual prompts. While effective, the web-crawled captions are usually short and incomplete, only describing some foreground information in the image. Thus, the VPG is only learned to extract obvious information for common captioning and might fail to extract specific information required by particular tasks. Furthermore, this issue becomes more severe in I4, where the tasks require the VPG to attend to specific visual information according to the other images in context (*e.g.,* express the subtle difference between two images).

To handle the critical issue of the VPG in existing MLLMs, we propose a lightweight Controllable Knowledge Re-Injection (CLORI) module that leverages the sophisticated reasoning ability of LLMs to control the VPG (*i.e.,* Q-former) to re-extract the missing visual information conditioned on instruction-specific semantics. Specifically, we first adopt the Q-former to generate task-agnostic visual prompts to provide basic information of the images for the LLM. Then, we derive instruction-specific conditions from the language model and take the condition to control the Q-former to conditionally extract specific information from the images, which is further re-injected into the LLM. Complemented by the conditional visual information, our model can effectively understand various intricate vision-language instructions. Our CLORI module is **computation-efficient** as it re-injects the newly extracted information at the same forward pass of the LLM, and is **training-efficient** as it reuses the frozen Q-former and only a set of conditional control parameters need to be fine-tuned, which are specially initialized to avoid noisy at the early training stage.

To methodically learn the CLORI module, we further introduce a Cross-Attention Guided Counterfactual Image Training (CAGIT) strategy, which automatically generates counterfactual training instances by collaborating a cascade of foundation models. Given an image, we first identify the most ignored areas by the Q-former based on its internal cross-attention maps. Then, we leverage ChatGPT (OpenAI, 2023a) and SAM (Kirillov et al., 2023) to determine the editing targets and generate a suitable type of editing description. Next, we adopt Blended Diffusion (Avrahami et al., 2022) to generate a counterfactual image by performing local edits on the original image based on the editing description. Finally, an inter-image discriminative pre-training task is formulated as describing the subtle difference between the original image and the generated counterfactual image. Since the edited parts are chosen from the most ignored areas, the CLORI module is forced to extract the missing visual information conditioned on the counterfactual image and the task instruction.

Empowered by controllable knowledge re-injection, we propose **Cheetor**🐆, a Transformer-based MLLM that can effectively compose holistic semantics from various intricate vision-language instructions. Without massive multi-modal instruction tuning data, the lightweight CLORI module can be effectively tuned by the CAGIT strategy along with image-text pairs of less than 1 million, which can be completed in several hours with a single A100 GPU. While computation- and data-efficient, our model significantly outperforms existing MLLMs on the complex I4 benchmark. Further, we evaluate Cheetor on the MME (Fu et al., 2023) benchmark, where our model demonstrates competitive performance.

Our contributions are summarized as follows: (1) We build **I4**, a comprehensive benchmark for interleaved vision-language instruction tuning of 31 tasks, covering a wide variety of real-world scenarios. (2) We propose a lightweight controllable knowledge re-injection (CLORI) module, which complementally re-injects instruction-specific visual information into the LLM following LLM-generated conditions. (3) We introduce a cross-attention guided counterfactual image training strategy, which can effectively learn the CLORI module using only 30k images. (4) Without high-quality multi-modal instruction tuning data, our Cheetor achieves state-of-the-art performance on the intricate I4 benchmark with the cost of 7 A100 GPU hours.

|          | Tasks | Scenarios | Images | Instructions | Avg. Images / Instruction | Avg. Words / Instruction |
|----------|-------|-----------|--------|--------------|---------------------------|--------------------------|
| **I4-Core** | 29 | 19 | 62.81K | 18.18K | 3.46 | 92.69 |
| **I4-Full** | 31 | 20 | 1.77M | 477.72K | 3.70 | 97.58 |

Table 1: Detailed statistics of I4 benchmark.

## 2 I4 BENCHMARK

To foster the research on general-purpose vision-language models, we manually collect a wide range of tasks that involve complex interleaved image-text sequential context and transform them into a unified instruction-response format.

**Data Format.** All task instances are given to the models in a unified instruction-response form to easily achieve zero-shot generalization on various tasks. Formally, each instance in I4 is composed of the following components:

- `Task_Instruction`: provides a complete natural language definition of a given task, including the input/output format and the task objective.

- `Task_Instance`: is a concrete sample of a given task that consists of interleaved image-text sequential context (*e.g.*, visually-rich textbooks and webpages, specific questions about the context).

- `Response`: represents the target output in natural language for a given task instruction and task instance. For classification tasks, we convert the class labels as options into the instruction and ask the model to output the option index in natural language as the response.

Without any specific emphasis, we use the term "instruction" to refer to the combination of `Task_Instruction` and `Task_Instance`. For each task, we manually design 10 `Task_Instruction` templates in natural language to increase the diversity.

**Task Collection and Categorization.** To comprehensively benchmark the interleaved vision-language instruction-following ability, we extensively gather a wide variety of multi-modal datasets from different fields and scenarios. As illustrated in Figure 1, our I4 benchmark covers 31 tasks of 7 categories across various scenarios (*i.e.*, surveillance, webpage, industrial, cartoon, *etc.*). Note that some datasets (*i.e.*, ALFRED, VISION, OCR-VQA) are not originally proposed for the task that involves interleaved image-text sequences. To further increase task diversity, we meticulously design certain rules to transform them to desired tasks, strictly following the original annotations.

**Evaluation Protocols.** Thanks to the unified task format of I4, all tasks can be evaluated in a zero-shot manner. For the open-ended generation tasks, we adopt *ROUGE-L* for evaluation. For the tasks that require the models to output option indexes, we take the *Accuracy* as the evaluation metric. While well-formated options are provided, we empirically observe that many MLLMs struggle to strictly follow instructions to output the option indexes but generate free-form text. Thus, when models do not exactly output the required options, we match their outputs to one of the given options based on the TF-IDF distance, which we find is more robust than model-based methods (OpenAI, 2023a; Reimers & Gurevych, 2019). Since we explore quantities of tasks, we take maximally 500 instances per task for evaluation efficiency and exclude several datasets that are difficult to obtain and are subject to strict copyright restrictions (referred as **I4-Core**). Meanwhile, we report the full version of the benchmark to facilitate future research on large-scale multi-modal instruction tuning (referred as **I4-Full**). Without special declaration, we use I4 to refer to I4-Core in the following.

**Benchmark Analysis.** Table 1 details the statistics for the benchmark. In total, I4-Full includes 31 tasks and 477.72K instruction-response pairs, serving as a large-scale benchmark for interleaved vision-language instruction following. On average, each instruction contains 3.70 images and 97.58 words. We report detailed information of the 31 tasks in Appendix.

Figure 2: Model architecture of the proposed Cheetor.

## 3 METHOD

### 3.1 OVERALL ARCHITECTURE

As illustrated in Figure 2, Cheetor is a vision-language instruction following model built upon the recent advanced language and vision-language foundation models. Specifically, we adopt the pre-trained ViT (Dosovitskiy et al., 2020) as our visual encoder, the pre-trained Q-Former from BLIP-2 (Li et al., 2023b) as our visual prompt generator, and the LLaMA2 (Touvron et al., 2023b) and Vicuna (Chiang et al., 2023) as our language decoder. The Q-former extracts relevant visual prompts from the dense ViT features as input to the language decoder, which serves as a general-purpose interface to unify various vision-language tasks as free-text generation.

To effectively comprehend instructions that involve intricate image-text sequential context, we propose a parameter-efficient controllable knowledge re-injection (CLORI) module (Section 3.2), which can additionally draw missing visual information based on instruction-specific conditions. Then, we introduce a cross-attention guided counterfactual image training (CAGIT) strategy (Section 3.3) to efficiently tune the CLORI module without the need of any instruction tuning data.

### 3.2 CONTROLLABLE KNOWLEDGE RE-INJECTION

Due to the pre-training objective of current VPGs (image-captioning alignment), they tend to extract significant foreground information as visual prompts for the language decoder, thus failing to provide specific visual information required by particular tasks. This phenomenon is even more severe in I4 benchmark as many tasks require the VPG to attend to specific information conditioned on visual cues from the other images. To address this limitation, InstructBLIP (Dai et al., 2023) further fine-tunes the Q-former to extract visual features according to instructions using 16M multi-modal instruction tuning data. While achieving outstanding performance on in-domain tasks, a recent study (Xu et al., 2023) indicates that fine-tuning on massive in-domain data severely undermines its generalizability on open-world scenarios.

**Overview.** Instead of directly relying on the Q-former to achieve task-specific feature extraction by massive instruction tuning, we introduce a lightweight CLORI module that utilizes the sophisticated reasoning ability of LLMs to control the Q-former to conditionally extract specific visual features, and further re-inject them into the LLM. As illustrated in Figure 2, we first take the Q-former to generate task-agnostic visual prompts to enable the LLM to form basic understanding of the given multi-modal instruction. Then, we derive instruction-specific conditions from the middle layer of the LLM, which is further used to control the Q-former to re-extract specific information from the images. Finally, the newly extracted visual information is injected into the LLM.

**Conditional Control Generation.** Given a multi-modal instruction, we first adopt the BLIP-2 pre-trained Q-former to generate general visual prompts for each image in the instruction. Note that, we use the Q-former without instruction data tuning as we aim to extract the task-agnostic foreground information at the first time. Q-former takes a fixed number of $K$ learnable queries to interact with image features by several cross-attention layers, and the output query representations are used as visual prompts, which are inserted into the position of their corresponding images in the instruction. We denote the input instruction for the language decoder as $H^0 = \{h_1^0, h_2^0, ..., v_{11}^0, ..., v_{1K}^0, ..., h_i^0, ..., v_{j1}^0, ..., v_{jK}^0, ..., h_N^0\}$, where $h_i^0$ represents the $i$-th text token and

Figure 3: Pipeline demonstration of cross-attention guided counterfactual image training strategy.

$V_j^0 = \{v_{j1}^0, ..., v_{jK}^0\}$ represents the visual prompts for the $j$-th interleaved image. Taking the instruction as input to the $L$-layer language decoder, we then extract the output hidden representation of the last token $h_N^{L/2}$ at the $\frac{L}{2}$-th layer, which has sufficiently reasoned over the whole multi-modal context during the first $\frac{L}{2}$ layers and contains comprehensive instruction-aware semantics. Next, we infer the instruction-specific condition $c$ from $h_N^{L/2}$ via a linear projection layer: $c = \mathbf{Linear}(h_N^{L/2})$.

**Controllable Visual Knowledge Extraction and Injection.** After obtaining the instruction-specific condition from the language decoder, we compose it with a set of learnable queries: $c + Q$, where $Q \in \mathbf{R}^{K \times d}$ and $c$ is added to each query of $Q$. Then, we reuse the same Q-former with the conditionally generated queries to re-extract specific visual information, thus obtaining the visual prompts $\overline{V}_j = \{\overline{v}_{j1}, ..., \overline{v}_{jK}\}$ for each image $j$, which contains the complementary information missed by the original visual prompts. Finally, we re-inject $\overline{V}_j$ into the language decoder by incorporating them into the hidden representations of corresponding visual prompts: $\tilde{V}_j^{L/2} = V_j^{L/2} + \mathbf{Linear}(\overline{V}_j)$, which is taken as the input to the $(\frac{L}{2} + 1)$-th layer.

**Efficient Training.** Our CLORI module is parameter-efficient as the Q-former is frozen and only a set of query embeddings and two linear projection layers need to be fine-tuned (6.3M). To stabilize the training process (Zhang & Agrawala, 2023), we initialize the linear projection layers with zeros. Thus, at the early training stage, the input to the $(\frac{L}{2}+1)$-th layer can be converted to: $\tilde{V}_j^{L/2} = V_j^{L/2}$, which will not cause any influence to the LLM decoder.

### 3.3 CROSS-ATTENTION GUIDED COUNTERFACTUAL IMAGE TRAINING

**Overview.** The CAGIT strategy diagnoses the initially ignored areas by Q-former according to the cross-attention maps between the queries and the image features, and generates a counterfactual image by performing several types of editing on the ignored areas. Then, a multi-image pre-training task is formulated as describing the subtle difference between the original image and the counterfactual image. Considering the edits are performed in the areas that are mostly ignored by Q-former, our CLORI module is forced to re-extract the missing information controlled by the instruction-specific conditions. An overview is illustrated in Figure 3.

**Editing Target Identification.** The Q-former takes the queries to interact with frozen image features through several cross-attention layers and uses the output query representations as the visual prompts. Therefore, the cross-attention maps between queries and image features reflect the interest of queries. We average the cross-attention maps across all layers and all queries to obtain the global cross-attention map $A$, where the value $A_{ij}$ indicates the significance degree of the corresponding image feature by the original task-agnostic Q-former queries. After that, we employ the advanced vision foundation models to obtain all the objects with segmentation masks in the image. Then, the significance degree of each object $\Phi(o_i)$ is computed based on the cross-attention map $A$ with RoIAlign (He et al., 2017), where we average the values of $A$ within the mask $m_i$ by interpolation. $\Phi(o_i)$ reflects what degree the visual features of object $o_i$ is extracted by the Q-former. Thus, we select the most ignored objects based on the $\Phi(o_i)$ value.

**Editing Description Generation.** We define four types of editing: *modifying objects, swapping objects, deleting objects, and adding objects.* Given the selected object, we instruct ChatGPT (OpenAI, 2023a) to generate a suitable editing description that is in harmony with the context, where ChatGPT is prompted with the corresponding image caption and detailed object information (*i.e.,* labels, positions). For adding objects, we only select `BACKGROUND` objects.

**Counterfactual Image Generation.** Having the editing instruction, we generate the counterfactual image using a text-to-image latent diffusion model (*i.e.,* Blended Diffusion (Avrahami et al., 2022)). Blended Diffusion performs local editing on the image according to the target object mask and generated editing description, thus rendering the counterfactual image. To ensure quality, we filter the edited images using CLIP similarity.

**Inter-Image Discriminative Pre-Training.** Given the original image and the counterfactual image pair, along with the task instruction (*"Describe the difference between the images"*), the inter-image discriminative pre-training task is defined as generating sentences to describe the subtle difference between the images. We convert the editing description to acquire the ground-truth sentences.

## 3.4 IMPLEMENTATION DETAILS

**Model.** We choose ViT-G/14 from EVA-CLIP (Fang et al., 2023) as our visual encoder and pre-trained Q-former from BLIP-2 without instruction tuning as the task-agnostic visual prompt generator. For the large language model, we implement two versions: LLaMA2-7B (Touvron et al., 2023b) and Vicuna-7B (Chiang et al., 2023), with 32 Transformer layers, respectively. We derive instruction-specific conditions from the 16th layer and re-inject the conditional visual knowledge into the 17th layer.

**Training.** We keep the visual backbone, visual prompt generator, and the language model frozen, and tune the CLORI module using the proposed CAGIT strategy. Since BLIP-2 models do not include pre-trained Q-former that matches Vicuna and LLaMA2, we reuse the Q-former that matches FlanT5-XXL and fine-tune the last linear projection layer with 5 million image-text pairs to align it with Vicuna/LLaMA2. All the tunable parameters of our CLORI module are a set of query embeddings and two linear projection layers, which only accounts for 0.09% ($\sim$6.3M) of the entire model. As for CAGIT strategy, we select about 30k images that contain significantly ignored objects and perform different types of editing on them. Finally, we generate approximately 64k counterfactual images with suitable modifications. Without high-quality multi-modal instruction data like Instruct-BLIP and mPLUG-owl, we only use 64k CAGIT data and 0.7 million image-text data to train the CLORI module. We tune the CLORI module for 18k steps using a batch size 24 for CAGIT and 64 for image-text data, which takes about 7 hours to complete with a single A100 GPU. Additionally, we adopt the AdamW optimizer with $\beta = (0.9, 0.999)$, and set the learning rate and weight decay to 0.00002 and 0.05, respectively. We warm up the training with 2k warm-up steps, followed by a learning rate decay mechanism with the cosine schedule.

## 4 EXPERIMENTS

### 4.1 ZERO-SHOT EVALUATION ON I4 BENCHMARK

**Comparison with Advanced MLLMs.** In this section, we conduct comprehensive evaluation of our Cheetor and the recent advanced MLLMs on the proposed I4 benchmark. For all methods, we choose versions with parameter sizes less than 10B. Please refer to Appendix for details. The average results of each task category are summarized in Table 2, which indicates the following.

- Our Cheetor consistently outperforms all previous work by a large margin across all categories of tasks, which demonstrates the stronger generalizability to follow such complicated interleaved vision-language instructions.

- While previous works mainly fine-tune on massive multi-modal instruction tuning data, our Cheetor still achieves new state-of-the-art performance by efficiently empowering the VPG using only thousands of images. This validates the effectiveness of the instruction-specific visual features conditionally extracted by our CLORI module, which provides complementary visual information for the language model to fully understand the complex sequential context of the instruction.

Table 2: Average results of zero-shot evaluation on each task category of I4 Benchmark.

| | Multi-Modal Dialogue | Visual Relation Inference | Visual Storytelling | Multi-Modal Cloze | Knowledge Grounded QA | Text-Rich Images QA | Multi-Image Reasoning |
|---|---|---|---|---|---|---|---|
| BLIP-2 | 11.96 | 20.10 | 3.67 | 18.25 | 39.73 | 30.53 | 39.53 |
| InstructBLIP | 33.58 | 24.41 | 11.48 | 21.20 | 47.40 | 44.40 | 48.55 |
| LLaMA-Adapter V2 | 14.22 | 17.57 | 13.51 | 18.00 | 44.80 | 32.00 | 44.03 |
| LLaVA | 7.79 | 10.70 | 8.27 | 15.85 | 36.20 | 28.33 | 41.53 |
| MiniGPT-4 | 13.69 | 17.07 | 7.95 | 16.60 | 30.27 | 26.40 | 43.50 |
| mPLUG-Owl | 12.67 | 19.33 | 5.40 | 16.25 | 33.27 | 32.47 | 42.50 |
| OpenFlamingo | 16.88 | 24.22 | 13.85 | 21.65 | 32.00 | 30.60 | 41.63 |
| Otter | 15.37 | 15.57 | 11.39 | 16.00 | 41.67 | 27.73 | 43.85 |
| **Cheetor-LLaMA2-7B** | **42.70** | 24.76 | 25.50 | **22.95** | **51.00** | **44.93** | 48.68 |
| **Cheetor-Vicuna-7B** | 37.50 | **25.20** | **25.90** | 22.15 | 48.60 | **44.93** | **50.28** |

- Compared with previous works that fine-tune the large-scale language decoder or visual encoder (*i.e.,* LLaVA, mPLUG-Owl), our model only tunes the lightweight CLORI module with 6.3M parameters and achieves significant performance gain.

- Cheetor exhibits significant superiority in several challenging tasks. For instance, Cheetor surpasses the SOTA methods by 3.6% on knowledge grounded QA, which requires models to infer answers from various multi-modal materials (*i.e.,* visually-rich webpages and textbooks).

**Innovative Findings.** The extensive evaluation on I4 benchmark reveals several key findings.

- **Limited Instruction Following Ability.** Despite existing vision-language models leveraging state-of-the-art LLMs, which have demonstrated impressive ability in following language instructions, this competence seems to falter when dealing with complex multi-modal instructions. For instance, when tasked with selecting the correct answer from a choice list given the context of images and texts, we observed some models inclining more towards describing the contents of the images instead of addressing the posed questions. This is perceived as a deficiency in the image-text alignment training process, to which we attribute the discrepancy.

- **Poor Performance on Interleaved Vision-Language Instructions.** While several models (*e.g.,* OpenFlamingo, Otter, mPLUG-owl) have been tuned on interleaved vision-language data, they still struggle to perform well on the interleaved vision-language instructions of the I4 benchmark. For instance, tasks such as describing the scene of a single image, often considered simple for current MLLMs, become increasingly challenging when it comes to the visual relation inference task. This particular task demands the model to compare scenes from two images and summarize their relationships. Since most MLLM's VPGs are learned to extract common foreground information for captioning, it might be insufficient for them to infer the relationships (*e.g.,* subtle differences between two images) as some details are ignored. Thus, the model's performance is impeded, leading to a decline in the efficacy of its multi-image reasoning ability.

- **Failing to Process Image-Choice Questions.** When dealing with multi-modal cloze tasks, we find that all models are limited to processing instructions that involve images as options. The performance is closed to random selection. We hope future work to utilize the new benchmark to make progress on this type of interleaved vision-language instructions.

## 4.2 ZERO-SHOT EVALUATION ON MME BENCHMARK

We evaluate our Cheetor on the recently proposed MME benchmark to further illustrate its strong generalizability to follow a diverse range of instructions. MME benchmark measures both perception and cognition abilities on a total of 14 subtasks. We report the averge results of perception tasks and cognition tasks in Table 3, respectively. While we do not fine-tune our model using massive multi-modal instruction tuning data, our Cheetor achieves competitive performance, compared with the intruction tuned models. Particularly, Cheetor exhibits superior performance on the perception tasks, which indicates the effectiveness of the proposed controllable knowledge re-injection mechnism. Please refer to Appendix for detailed results.

Table 3: Zero-shot evaluation of perception and cognition abilities on MME benchmark.

| | BLIP-2 | InstructBLIP | LA-V2 | LLaVA | MiniGPT-4 | mPLUG-Owl | Otter | **Cheetor** |
|---|---|---|---|---|---|---|---|---|
| Perception | 1293.84 | 1212.82 | 972.67 | 502.82 | 866.57 | 967.34 | 1292.26 | 1299.24 |
| Cognition | 290.00 | 291.79 | 248.93 | 214.64 | 292.14 | 276.07 | 306.43 | 321.07 |

Figure 4: Qualitative examples generated by our Cheetor-Vicuna-7B model.

## 4.3 QUALITATIVE EVALUATION

As illustrated in Figure 4, our Cheetor demonstrates strong abilities to perform reasoning over complicated interleaved vision-language instructions. For instance, in **(a)**, Cheetor is able to keenly identify the connections between the images and thereby infer the reason that causes this unusual phenomenon. In **(b, c)**, Cheetor can reasonably infer the relations among the images and understand the metaphorical implications they want to convey. In **(e, f)**, Cheetor exhibits the ability to comprehend absurd objects through multi-modal conversations with humans.

## 5    RELATED WORK

MLLMs (Yin et al., 2023) aim to serve as a general-purpose assistant to perform various vision-language tasks by free-text generation. Flamingo (Alayrac et al., 2022) and BLIP-2 (Li et al., 2023b) bridge LLMs with powerful pre-trained visual encoders and demonstrate strong zero-shot ability by aligning visual features with LLMs. Inspired by the great success of instruction-tuned LLMs, follow-up works of LLaVA (Liu et al., 2023), MiniGPT-4 (Zhu et al., 2023), InstructBLIP (Dai et al., 2023), Otter (Li et al., 2023a), mPLUG-Owl (Ye et al., 2023), OpenFlamingo (Awadalla et al., 2023) propose to fine-tune MLLMs with multi-modal instruction tuning data. To effectively benchmark the recent progress in MLLMs, concurrent works of LVLM-eHub (Xu et al., 2023) and MME Benchmark (Fu et al., 2023) are proposed, while they mainly focus on vision-language instruction data that only involves a single image with limited instruction diversity. In this paper, we propose the first interleaved vision-language instruction-following benchmark with three characteristics: *interleaved vision-language context, diverse forms of complex instructions,* and *a vast range of instruction-following scenarios.* Furthermore, we propose a lightweight controllable knowledge re-injection module to address the inherent limitation of current VPGs. Our parameter-efficient controllable knowledge re-injection module is efficiently tuned by our proposed cross-attention guided counterfactual image training strategy, which demonstates powerful potentials of text-to-image diffusion models (He et al., 2022; Lin et al., 2023; Prabhu et al., 2023; Bansal & Grover, 2023; Yu et al., 2023) to facilitate vision-language understanding (Radford et al., 2021b; Li et al., 2022a).

## 6    CONCLUSION

In this paper, we introduce **I4**, a comprehensive evaluation benchmark for multimodal large language models, consisting of 31 tasks with complicated vision-language sequential context, covering a wide range of scenarios. After systematically evaluating 7 advanced MLLMs, we propose a parameter-efficient controllable knowledge re-injection module with an annotation-free cross-attention guided counterfactual image training strategy to address a common defect in their visual prompt generators. Empowered by the proposed approach, the language model can selectively control the VPG to re-extract specific visual information conditioned on instruction semantics. The resulting **Cheetor** model achieves state-of-the-art instruction following ability on I4, without high-quality multi-modal instruction tuning data. Further, evaluation on concurrent MME benchmark shows that our Cheetor also exhibits competitive performance compared with state-of-the-art instruction tuned models.

REFERENCES

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18208–18218, 2022.

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`.

Haoping Bai, Shancong Mou, Tatiana Likhomanenko, Ramazan Gokberk Cinbis, Oncel Tuzel, Ping Huang, Jiulong Shan, Jianjun Shi, and Meng Cao. Vision datasets: A benchmark for vision-based industrial inspection. *arXiv preprint arXiv:2306.07890*, 2023.

Ankan Bansal, Yuting Zhang, and Rama Chellappa. Visual question answering on image sets. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*, pp. 51–67. Springer, 2020.

Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets. *arXiv preprint arXiv:2302.02503*, 2023.

Nilavra Bhattacharya, Qing Li, and Danna Gurari. Why does a visual question have different answers? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4271–4280, 2019.

Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.

Yingshan Chang, Mridu Narang, Hisami Suzuki, Guihong Cao, Jianfeng Gao, and Yonatan Bisk. Webqa: Multihop and multimodal qa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16495–16504, 2022.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90% chatgpt quality, 2023. URL `https://vicuna.lmsys.org`.

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.

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.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

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*, pp. 19358–19369, 2023.

Maxwell Forbes, Christine Kaeser-Chen, Piyush Sharma, and Serge Belongie. Neural naturalist: generating fine-grained image comparisons. *arXiv preprint arXiv:1909.04101*, 2019.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.

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.

Tanmay Gupta, Dustin Schwenk, Ali Farhadi, Derek Hoiem, and Aniruddha Kembhavi. Imagine this! scripts to compositions to videos. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 598–613, 2018.

Xintong Han, Zuxuan Wu, Phoenix X Huang, Xiao Zhang, Menglong Zhu, Yuan Li, Yang Zhao, and Larry S Davis. Automatic spatially-aware fashion concept discovery. In *Proceedings of the IEEE international conference on computer vision*, pp. 1463–1471, 2017.

Darryl Hannan, Akshay Jain, and Mohit Bansal. Manymodalqa: Modality disambiguation and qa over diverse inputs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 7879–7886, 2020.

Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.

Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. Is synthetic data from generative models ready for image recognition? *arXiv preprint arXiv:2210.07574*, 2022.

Mehrdad Hosseinzadeh and Yang Wang. Image change captioning by learning from an auxiliary task. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2725–2734, 2021.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pp. 1233–1239, 2016.

Phillip Isola, Joseph J Lim, and Edward H Adelson. Discovering states and transformations in image collections. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1383–1391, 2015.

Mohit Iyyer, Varun Manjunatha, Anupam Guha, Yogarshi Vyas, Jordan Boyd-Graber, Hal Daume, and Larry S Davis. The amazing mysteries of the gutter: Drawing inferences between panels in comic book narratives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern recognition*, pp. 7186–7195, 2017.

Harsh Jhamtani and Taylor Berg-Kirkpatrick. Learning to describe differences between pairs of similar images. *arXiv preprint arXiv:1808.10584*, 2018.

Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern recognition*, pp. 4999–5007, 2017.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*, 2023a.

Juncheng Li, Xin He, Longhui Wei, Long Qian, Linchao Zhu, Lingxi Xie, Yueting Zhuang, Qi Tian, and Siliang Tang. Fine-grained semantically aligned vision-language pre-training. *Advances in neural information processing systems*, 35:7290–7303, 2022a.

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*, 2023b.

Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6329–6338, 2019.

Yongqi Li, Wenjie Li, and Liqiang Nie. Mmcoqa: Conversational question answering over text, tables, and images. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4220–4231, 2022b.

Shaobo Lin, Kun Wang, Xingyu Zeng, and Rui Zhao. Explore the power of synthetic data on few-shot object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 638–647, 2023.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.

Adyasha Maharana, Darryl Hannan, and Mohit Bansal. Storydall-e: Adapting pretrained text-to-image transformers for story continuation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*, pp. 70–87. Springer, 2022.

Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2200–2209, 2021.

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)*, pp. 947–952. IEEE, 2019.

OpenAI. Chatgpt: A language model for conversational ai. Technical report, OpenAI, 2023a. URL https://www.openai.com/research/chatgpt.

OpenAI. Gpt-4 technical report. *arXiv:2303.08774*, 2023b.

Viraj Prabhu, Sriram Yenamandra, Prithvijit Chattopadhyay, and Judy Hoffman. Lance: Stress-testing visual models by generating language-guided counterfactual images. *arXiv preprint arXiv:2305.19164*, 2023.

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*, pp. 8748–8763. PMLR, 2021a.

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*, pp. 8748–8763. PMLR, 2021b.

Hareesh Ravi, Kushal Kafle, Scott Cohen, Jonathan Brandt, and Mubbasir Kapadia. Aesop: Abstract encoding of stories, objects, and pictures. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2052–2063, 2021.

Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992, 2019.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10740–10749, 2020.

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.

Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. Multimodalqa: Complex question answering over text, tables and images. *arXiv preprint arXiv:2104.06039*, 2021.

Hao Tan, Franck Dernoncourt, Zhe Lin, Trung Bui, and Mohit Bansal. Expressing visual relationships via language. *arXiv preprint arXiv:1906.07689*, 2019.

Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku Hasegawa, Itsumi Saito, and Kuniko Saito. Slidevqa: A dataset for document visual question answering on multiple images. *arXiv preprint arXiv:2301.04883*, 2023.

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:2302.13971*, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023b.

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.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.

Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9068–9079, 2018.

Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*, 2023.

Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. *arXiv preprint arXiv:1809.00812*, 2018.

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.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.

Qifan Yu, Juncheng Li, Wentao Ye, Siliang Tang, and Yueting Zhuang. Interactive data synthesis for systematic vision adaptation via llms-aigcs collaboration. *arXiv preprint arXiv:2305.12799*, 2023.

Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

## A  Benchmark Details

| Task | Scenario | Dataset | Metirc |
|------|----------|---------|--------|
| **Multi-Modal Dialogue** | | | |
| Conversational Embodied Dialogue | Embodied | ALFRED (Shridhar et al., 2020) | ROUGE-L |
| Multi-Modal Dialogue | Conversation | MMCoQA (Li et al., 2022b) | ROUGE-L |
| **Visual Relation Inference** | | | |
| Visual Change Captioning | Surveillance | Spot-the-Diff (Jhamtani & Berg-Kirkpatrick, 2018) | ROUGE-L |
| Visual Change Captioning | Synthetic | CLEVR-Change (Hosseinzadeh & Wang, 2021) | ROUGE-L |
| Visual Relationship Expressing | General | IEdit (Tan et al., 2019) | ROUGE-L |
| Subtle Difference Expressing | Fine-Grained | Birds-to-Words (Forbes et al., 2019) | ROUGE-L |
| **Visual Storytelling** | | | |
| Animated Story Completion | Cartoon | AESOP (Ravi et al., 2021) | ROUGE-L |
| Animated Story Completion | Cartoon | PororoSV (Li et al., 2019) | ROUGE-L |
| Animated Story Completion | Cartoon | FlintstonesSV (Gupta et al., 2018) | ROUGE-L |
| Sequential Photo Storytelling | Album | VIST (Huang et al., 2016) | ROUGE-L |
| Sequential Photo Storytelling | Cartoon | DiDeMoSV (Maharana et al., 2022) | ROUGE-L |
| **Multi-Modal Cloze** | | | |
| Comic Dialogue Identification | Cartoon | COMICS-Dialogue (Iyyer et al., 2017) | Accuracy |
| Comic Panel Identification | Cartoon | COMICS-Panel (Iyyer et al., 2017) | Accuracy |
| Recipe Completion | Recipe | RecipeQA-TextCloze (Yagcioglu et al., 2018) | Accuracy |
| Visual Step Cloze | Recipe | RecipeQA-VisualCloze (Yagcioglu et al., 2018) | Accuracy |
| **Knowledge Grounded QA** | | | |
| Webpage QA | Webpage | WebQA (Chang et al., 2022) | Accuracy |
| Textbook QA | Textbook | TQA (Kembhavi et al., 2017) | Accuracy |
| Complex Multimodal QA | Wikipedia | MMQA (Talmor et al., 2021) | Accuracy |
| Complex Multimodal QA* | Wikipedia | MANYMODALQA (Hannan et al., 2020) | Accuracy |
| **Text-Rich Images QA** | | | |
| Slide QA | Slide | SlideVQA (Tanaka et al., 2023) | Accuracy |
| OCR QA | Book Cover | OCR-VQA (Mishra et al., 2019) | Accuracy |
| Document QA | Document Image | DocVQA (Mathew et al., 2021) | Accuracy |
| **Multi-Image Reasoning** | | | |
| Image-Set QA* | Indoor Egocentric | Gibson (Bansal et al., 2020; Xia et al., 2018) | Accuracy |
| Image-Set QA | Driving Recording | nuScenes (Bansal et al., 2020; Caesar et al., 2020) | Accuracy |
| Industrial Inspection | Industrial | VISION (Bai et al., 2023) | Accuracy |
| Fashion QA | Fashion | Fashion200K (Han et al., 2017) | Accuracy |
| Property Coherence | General | MIT-States-PropertyCoherence (Isola et al., 2015) | Accuracy |
| State Transformation Coherence | General | MIT-States-StateCoherence (Isola et al., 2015) | Accuracy |
| Visual Step Matching | Recipe | RecipeQA-ImageCoherence (Yagcioglu et al., 2018) | Accuracy |
| Multi-Image Visual Entailment | General | NLVR2 (Suhr et al., 2018) | Accuracy |
| Ambiguity Analysis | Mobile Photo | VizWiz (Bhattacharya et al., 2019) | Accuracy |

Table 4: Summary of the interleaved vision-language instruction-following tasks in I4 benchmark. * indicates the tasks that are not included in I4-Core.

## B  Model Details in I4 Benchmark

- **LLaVA** (Liu et al., 2023) establishes a connection between the visual encoder ViT-L/14 from CLIP (Radford et al., 2021a) and the language decoder LLaMA (Touvron et al., 2023a), utilizing a lightweight, fully-connected (FC) layer. Initially, the system trains this FC layer using 595K image-text pairs, while keeping both the visual encoder and LLM static. Following this, LLaVA fine-tunes both the FC layer and LLM using a dataset comprising 158K instructional vision-language pairs. The tested version is "LLaVA-7B-v0".

- **LLaMA-Adapter V2** (Gao et al., 2023) stands as a model of parameter efficiency within the realm of visual instruction. Despite maintaining the visual encoder (ViT-L/14) and the LLM in a static state, LA-V2 distributes the instruction-following capacity of the entire LLaMA system via bias-tuning. This method allows for the refinement of scale, bias, norm, and prompt parameters on diverse data sets. These include 200M image captioning data, 158K visual instruction-following data, and an additional 52K language instruction-following data, the latter of which was assembled by GPT-4 (OpenAI, 2023b). The tested version is "LLaVA-7B".

- **MiniGPT-4** (Zhu et al., 2023) bridges the gap between the visual encoder and text encoder using a fully-connected (FC) layer. Initially, this model trains the FC layer on a

dataset comprised of 5M image-text pairs before fine-tuning it on 3.5K instructional vision-language data. Notwithstanding its simplicity, MiniGPT-4 requires the loading of a pre-trained vision encoder from BLIP2, as well as a Vicuna LLM (Chiang et al., 2023). The tested version is "minigpt4-aligned-with-vicuna7b".

- **BLIP2** (Li et al., 2023b) employs a dual-stage strategy to seamlessly bridge the modality gap, utilizing a lean Q-Former pre-trained on 129 million image-text pairs. The initial stage kick-starts the learning process of vision-language representation, leveraging a frozen image encoder, the ViT-g/14 from EVA-CLIP (Fang et al., 2023). Subsequently, the second stage harnesses a frozen LLM, the Vicuna-7B (Chung et al., 2022), to initiate the vision-to-language generative learning. This innovative strategy effectively facilitates zero-shot instructed image-to-text generation. The tested version is "blip2-pretrained-vicuna7b".

- **mPLUG-Owl** (Ye et al., 2023) introduces a visual abstractor, fundamentally close the Perceiver Resampler in Flamingo (Alayrac et al., 2022), as a bridge between the pre-trained visual encoder ViT-L/14 and the LLM (LLaMA (Touvron et al., 2023a)). This model adopts a two-stage fine-tuning procedure. In the initial phase, both the visual encoder and the visual abstractor undergo comprehensive fine-tuning using a dataset of 204M image-text pairs. Subsequently, in the second phase, mPLUG-Owl applies the 158K LLaVA-Instruct dataset to fine-tune the pre-trained LLM in a parameter-efficient manner through the use of LoRA (Hu et al., 2021). The tested version is "mplug-owl-llama-7b".

- **Otter** (Li et al., 2023a) is a multimodal model that applies in-context instruction tuning based on OpenFlamingo (Alayrac et al., 2022). This model integrates a LLaMA-7B (Touvron et al., 2023a) language encoder and a CLIP ViT-L/14. While the visual and text encoders remain static, Otter refines an additional 1.3 billion parameters. These parameters are derived from adaptation modules and are trained using 158K instruction-following data. The tested version is "OTTER-Image-LLaMA7B-LA-InContext".

- **InstructBLIP** (Dai et al., 2023) originates from a pre-trained BLIP-2 model, which consists of a ViT-g/14 image encoder, a Vicuna LLM, and a Q-Former to act as the bridge between these two components. During the process of vision-language instruction tuning, only the Q-Former undergoes fine-tuning, with the training process leveraging data from 13 distinct visual question-answering datasets. The tested version is "blip2-vicuna-instruct-7b".

- **OpenFlamingo** (Alayrac et al., 2022; Awadalla et al., 2023) represents one of the pioneering efforts to incorporate Language Model Learning (LLMs) into the domain of vision-language pretraining. To optimize its conditioning on visual features, Flamingo strategically integrates a number of gated cross-attention dense blocks amidst the layers of the pre-trained language encoder. OpenFlamingo offers an open-source rendition of this advanced model. The tested version is "llama-7b".

## C  DETAILED ZERO-SHOT PERFORMANCE ON MME BENCHMARK

Table 5: Detailed zero-shot performance on MME benchmark.

|  | BLIP-2 | InstructBLIP | LA-V2 | LLaVA | MiniGPT-4 | mPLUG-Owl | Otter | **Cheetor** |
|---|---|---|---|---|---|---|---|---|
| Existence | 160.00 | 185.00 | 120.00 | 50.00 | 115.00 | 120.00 | 195.00 | 180.00 |
| Count | 135.00 | 143.33 | 50.00 | 50.00 | 123.33 | 88.33 | 50.00 | 96.67 |
| Position | 73.33 | 66.67 | 48.33 | 50.00 | 81.67 | 50.00 | 86.67 | 80.00 |
| Color | 148.33 | 153.33 | 75.00 | 55.00 | 110.00 | 55.00 | 113.33 | 116.67 |
| Poster | 141.84 | 123.81 | 99.66 | 50.00 | 55.78 | 136.05 | 138.78 | 147.28 |
| Celebrity | 105.59 | 101.18 | 86.18 | 48.82 | 65.29 | 100.29 | 172.65 | 164.12 |
| Scene | 145.25 | 153.00 | 148.50 | 50.00 | 95.75 | 135.50 | 158.75 | 156.00 |
| Landmark | 138.00 | 79.75 | 150.25 | 50.00 | 69.00 | 159.25 | 137.25 | 145.00 |
| Artwork | 136.50 | 134.25 | 69.75 | 49.00 | 55.75 | 96.25 | 129.00 | 113.50 |
| OCR | 110.00 | 72.50 | 125.00 | 50.00 | 95.00 | 65.00 | 72.50 | 100.00 |
| Perception | 1293.84 | 1212.82 | 972.67 | 502.82 | 866.57 | 967.34 | 1292.26 | 1299.24 |
| Commonsense | 110.00 | 129.29 | 81.43 | 57.14 | 72.14 | 78.57 | 106.43 | 98.57 |
| Numerical | 40.00 | 40.00 | 62.50 | 50.00 | 55.00 | 60.00 | 72.50 | 77.50 |
| Text Translation | 65.00 | 65.00 | 50.00 | 57.50 | 55.00 | 80.00 | 57.50 | 57.50 |
| Code Reasoning | 75.00 | 57.50 | 55.00 | 50.00 | 110.00 | 57.50 | 70.00 | 87.50 |
| Cognition | 290.00 | 291.79 | 248.93 | 214.64 | 292.14 | 276.07 | 306.43 | 321.07 |

# D DETAILED ZERO-SHOT PERFORMANCE ON I4 BENCHMARK

Table 6: Zero-shot evaluation on multi-modal dialogue.

|  | Conversational Embodied Dialogue | Multi-Modal Dialogue |
|---|---|---|
| BLIP-2 | 6.52 | 17.39 |
| InstructBLIP | 18.07 | 49.09 |
| LLaMA-Adapter V2 | 19.04 | 9.40 |
| LLaVA | 10.19 | 5.39 |
| MiniGPT-4 | 16.82 | 10.57 |
| mPLUG-Owl | 11.07 | 14.27 |
| OpenFlamingo | 24.27 | 9.49 |
| Otter | 16.06 | 14.68 |
| **Cheetor-LLaMA2-7B** | 48.31 | 37.04 |
| **Cheetor-Vicuna-7B** | 41.02 | 33.99 |

Table 7: Zero-shot evaluation on visual relation inference.

|  | Visual Change Captioning -Spot-the-Diff | Visual Change Captioning -CLEVR-Change | Visual Relationship Expressing | Subtle Difference Expressing |
|---|---|---|---|---|
| BLIP-2 | 0.40 | 0.08 | 9.27 | 4.89 |
| InstructBLIP | 19.71 | 4.61 | 10.70 | 10.92 |
| LLaMA-Adapter V2 | 16.72 | 15.52 | 7.88 | 13.92 |
| LLaVA | 8.50 | 8.76 | 6.72 | 9.11 |
| MiniGPT-4 | 7.50 | 7.49 | 7.84 | 8.97 |
| mPLUG-Owl | 6.06 | 1.46 | 6.22 | 7.86 |
| OpenFlamingo | 13.01 | 11.90 | 12.57 | 17.90 |
| Otter | 12.69 | 11.63 | 8.85 | 12.38 |
| **Cheetor-LLaMA2-7B** | 21.02 | 42.05 | 14.10 | 24.81 |
| **Cheetor-Vicuna-7B** | 20.01 | 41.60 | 16.35 | 25.64 |

Table 8: Zero-shot evaluation on visual storytelling.

|  | Animated Story Completion-AESOP | Animated Story Completion-PororoSV | Animated Story Completion-FlintstonesSV | Sequential Photo Storytelling-VIST | Sequential Photo Storytelling-DiDeMoSV |
|---|---|---|---|---|---|
| BLIP-2 | 22.64 | 25.04 | 28.61 | 12.22 | 11.98 |
| InstructBLIP | 18.80 | 28.20 | 33.32 | 16.92 | 24.80 |
| LLaMA-Adapter V2 | 18.01 | 20.15 | 24.22 | 10.89 | 14.57 |
| LLaVA | 13.56 | 11.44 | 12.77 | 8.00 | 7.71 |
| MiniGPT-4 | 12.23 | 16.00 | 26.48 | 14.82 | 15.81 |
| mPLUG-Owl | 18.28 | 20.49 | 32.12 | 10.82 | 14.94 |
| OpenFlamingo | 23.32 | 32.35 | 37.79 | 15.14 | 12.50 |
| Otter | 13.94 | 17.52 | 22.21 | 9.96 | 14.23 |
| **Cheetor-LLaMA2-7B** | 19.98 | 28.67 | 38.14 | 16.95 | 20.05 |
| **Cheetor-Vicuna-7B** | 19.93 | 28.36 | 39.19 | 17.34 | 21.27 |

Table 9: Zero-shot evaluation on multi-modal cloze.

| | Comic Dialogue Identification | Comic Panel Identification[1] | Recipe Completion | Visual Step Cloze[1] |
|---|---|---|---|---|
| BLIP-2 | 40.80 | 0.00 | 31.00 | 1.20 |
| InstructBLIP | 40.60 | 0.00 | 27.40 | 16.80 |
| LLaMA-Adapter V2 | 24.40 | 0.40 | 38.20 | 9.00 |
| LLaVA | 30.60 | 0.00 | 32.80 | 0.00 |
| MiniGPT-4 | 33.00 | 1.00 | 31.60 | 0.80 |
| mPLUG-Owl | 36.60 | 0.00 | 27.60 | 0.80 |
| OpenFlamingo | 38.40 | 1.20 | 29.00 | 18.00 |
| Otter | 29.00 | 0.00 | 35.00 | 0.00 |
| **Cheetor-LLaMA2-7B** | 36.80 | 1.80 | 51.80 | 1.40 |
| **Cheetor-Vicuna-7B** | 39.20 | 3.60 | 30.40 | 15.40 |

[1] For tasks with images as options, only responses that begin with the correct answer will be evaluated as correct.

Table 10: Zero-shot evaluation on knowledge grounded QA.

| | Webpage QA | Textbook QA | Complex Multimodal QA |
|---|---|---|---|
| BLIP-2 | 48.80 | 29.60 | 40.80 |
| InstructBLIP | 45.20 | 30.20 | 66.80 |
| LLaMA-Adapter V2 | 44.60 | 46.00 | 43.80 |
| LLaVA | 39.40 | 39.60 | 29.60 |
| MiniGPT-4 | 27.40 | 28.60 | 34.80 |
| mPLUG-Owl | 34.20 | 30.00 | 35.60 |
| OpenFlamingo | 37.80 | 32.40 | 25.80 |
| Otter | 45.00 | 39.00 | 41.00 |
| **Cheetor-LLaMA2-7B** | 49.40 | 42.40 | 61.20 |
| **Cheetor-Vicuna-7B** | 50.00 | 33.40 | 62.40 |

Table 11: Zero-shot evaluation on text-rich images QA.

| | Slide QA | OCR QA | Document QA |
|---|---|---|---|
| BLIP-2 | 43.00 | 2.00 | 46.60 |
| InstructBLIP | 42.00 | 44.20 | 47.00 |
| LLaMA-Adapter V2 | 43.00 | 3.40 | 49.60 |
| LLaVA | 38.80 | 2.60 | 43.60 |
| MiniGPT-4 | 35.20 | 7.20 | 36.80 |
| mPLUG-Owl | 35.60 | 22.60 | 39.20 |
| OpenFlamingo | 35.60 | 3.80 | 52.40 |
| Otter | 38.40 | 2.20 | 42.60 |
| **Cheetor-LLaMA2-7B** | 45.80 | 39.60 | 49.40 |
| **Cheetor-Vicuna-7B** | 46.80 | 39.40 | 48.60 |

Table 12: Zero-shot evaluation on multi-image reasoning.

| | Image-Set QA | Industrial Inspection | Fashion QA | Property Coherence | State Transformation Coherence | Visual Step Matching[1] | Multi-Image Visual Entailment | Ambiguity Analysis |
|---|---|---|---|---|---|---|---|---|
| BLIP-2 | 32.80 | 44.60 | 42.20 | 59.00 | 38.20 | 0.20 | 56.40 | 42.80 |
| instructblip7b | 65.00 | 50.60 | 44.40 | 59.20 | 59.40 | 11.60 | 55.20 | 43.00 |
| LLaMA-Adapter V2 | 41.60 | 55.00 | 45.60 | 48.80 | 63.00 | 0.00 | 54.80 | 43.40 |
| LLaVA | 29.60 | 53.00 | 45.20 | 50.40 | 59.20 | 0.80 | 50.80 | 43.20 |
| MiniGPT-4 | 30.40 | 59.80 | 49.20 | 52.00 | 57.80 | 0.20 | 50.60 | 48.00 |
| mPLUG-Owl | 29.20 | 54.20 | 45.80 | 50.00 | 60.60 | 0.00 | 55.00 | 45.20 |
| OpenFlamingo | 25.80 | 52.20 | 44.20 | 59.60 | 51.40 | 2.20 | 53.60 | 44.00 |
| Otter | 44.80 | 69.80 | 47.00 | 51.40 | 46.40 | 0.00 | 49.00 | 42.40 |
| **Cheetor-LLaMA2-7B** | 62.60 | 61.40 | 46.00 | 56.60 | 57.80 | 0.00 | 53.80 | 51.20 |
| **Cheetor-Vicuna-7B** | 67.20 | 48.80 | 50.00 | 60.80 | 60.00 | 0.20 | 57.80 | 57.40 |

[1] For tasks with images as options, only responses that begin with the correct answer will be evaluated as correct.