

ImageRef-VL: Enabling Contextual Image Referencing in Vision-Language Models

Jingwei Yi^{1*}, Junhao Yin², Ju Xu², Peng Bao³, Yongliang Wang², Wei Fan⁴, Hao Wang^{5†}

¹University of Science and Technology of China ²ByteDance

³Peking University ⁴University of Oxford ⁵Tsinghua University

yjw1029@mail.ustc.edu.cn weifan.oxford@gmail.com hao-wang20@mails.tsinghua.edu.cn
{yinjunhao,yufeng.1016,baopeng.peter,yongliang.wyl}@bytedance.com

Abstract

Vision-Language Models (VLMs) have demonstrated remarkable capabilities in understanding multimodal inputs and have been widely integrated into Retrieval-Augmented Generation (RAG) based conversational systems. While current VLM-powered chatbots can provide textual source references in their responses, they exhibit significant limitations in referencing contextually relevant images during conversations. In this paper, we introduce *Contextual Image Reference* – the ability to appropriately reference relevant images from retrieval documents based on conversation context – and systematically investigate VLMs’ capability in this aspect. We conduct the first evaluation for contextual image referencing, comprising a dedicated testing dataset and evaluation metrics. Furthermore, we propose ImageRef-VL, a method that significantly enhances open-source VLMs’ image referencing capabilities through instruction fine-tuning on a large-scale, manually curated multimodal conversation dataset. Experimental results demonstrate that ImageRef-VL not only outperforms proprietary models but also achieves an 88% performance improvement over state-of-the-art open-source VLMs in contextual image referencing tasks. Our code is available at <https://github.com/bytedance/ImageRef-VL>.

1 Introduction

In recent years, Vision-Language Models (VLMs) have achieved remarkable progress, enabling advanced multi-modal reasoning and generation from combined text and image inputs. Both close-source models, such as GPT-4o (Hurst et al., 2024) and Claude (Anthropic, 2024), and open-source models,

*This work was done when the author Jingwei Yi was at ByteDance Group for intern.

†Corresponding authors.

¹The response was generated with reference to <https://en.wikipedia.org/wiki/Brachiosaurus>.

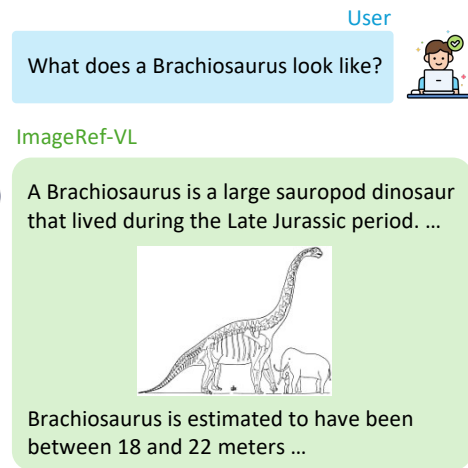


Figure 1: An example of contextual image reference, where referencing the images of a Brachiosaurus can largely enhance user comprehension and engagement.¹

such as Phi-3.5-Vision (Abdin et al., 2024), Qwen2-VL (Wang et al., 2024) and InternVL2 (Chen et al., 2024), have demonstrated impressive capabilities across a range of vision-language tasks (Hudson and Manning, 2019; Yu et al., 2023; Fu et al., 2024). Concurrently, the integration of VLMs with Retrieval-Augmented Generation (RAG) has allowed these models to retrieve and incorporate external knowledge into their responses²³⁴⁵. One feature of VLM-based RAG chatbots is to provide references to the retrieved text sources used to generate responses. Extensive research efforts have been dedicated to improving textual reference accuracy in RAG chatbots (Gao et al., 2023a; Shen et al., 2024; Zhang et al., 2024a; Fierro et al., 2024).

While VLM-based RAG systems have made significant strides in providing reliable textual references, they face a critical limitation in their ability

²<https://chatgpt.com/>

³<https://claude.ai/>

⁴<https://www.doubao.com/chat/>

⁵<https://www.perplexity.ai/>

to leverage visual content effectively during conversations. We identify this gap as the absence of *Contextual Image Reference* - the capability to strategically select and incorporate relevant images from retrieved documents to enhance response comprehension and user engagement. As demonstrated in Figure 1, when discussing complex subjects like the Brachiosaurus, purely textual descriptions of its physical characteristics often fail to convey information intuitively. Despite its potential impact on multimodal conversation systems, the challenge of contextual image referencing remains largely unexplored in current research.

To address these challenges, we first propose and formally define *Contextual Image Reference* as a novel task that requires VLMs to incorporate relevant images as integral components of their responses. To systematically evaluate performance on this task, we construct a dedicated testing dataset and develop comprehensive metrics that assess a model’s capability to integrate images in contextually appropriate ways. Building upon open-source VLMs, (i.e., InternVL2), we present ImageRef-VL, a framework that enhances models’ contextual image referencing abilities. Our approach involves generating initial responses using existing LLMs and VLMs, carefully curating these outputs through manual refinement, and leveraging the resulting high-quality dataset for supervised fine-tuning. Through this process, ImageRef-VL learns to make informed decisions about when and how to incorporate images as authoritative visual references.

The primary contributions of our work are summarized as follows:

- To the best of our knowledge, we are the first to introduce and formally define *Contextual Image Reference* as a novel task for multimodal conversational AI, addressing a critical gap in current VLM capabilities.
- We conduct a comprehensive evaluation for this task, including a carefully curated testing dataset and novel metrics that capture both the relevance and naturalness of image references.
- We propose ImageRef-VL, a fine-tuning framework that significantly advances the state-of-the-art in contextual image referencing, demonstrating an 88% performance improvement over existing open-source VLMs.

- Through extensive experiments across various scenarios, we validate the effectiveness of our approach and establish strong baseline results for future research.

2 Related Works

2.1 Vision Language Models

Large Language Models (LLMs), primarily based on the transformer architecture (Vaswani, 2017), have recently achieved remarkable performance in a wide range of natural language tasks (Zhang et al., 2023; Poesia et al., 2022; Kojima et al., 2022). Building upon these advances, Vision-Language Models (VLMs) extend LLM capabilities to the visual domain, enabling sophisticated reasoning and content generation from both textual and visual inputs (Meta, 2024; Lu et al., 2024; Liu et al., 2024b,a). Recently, numerous VLMs have been introduced, spanning both close-sourced systems (e.g., GPT-4o (Hurst et al., 2024), Claude (Anthropic, 2024)) and open-source frameworks (e.g., Phi-3.5-Vision (Abdin et al., 2024), Qwen2-VL (Wang et al., 2024), InternVL2 (Chen et al., 2024)). These models have demonstrated impressive capabilities across diverse vision-language benchmarks (Hudson and Manning, 2019; Yu et al., 2023; Fu et al., 2024), promoting an evolving research landscape in multimodal machine intelligence. Existing VLMs are generally composed of three core components, i.e., a vision encoder, an adapter, and an LLM backend. The vision encoder, such as CLIP (Radford et al., 2021) or BLIP (Li et al., 2022, 2023), is designed to extract rich visual features from images, converting them into representations that can be effectively processed by downstream components. The adapter component subsequently bridges these extracted features to the language model, employing architectures such as simple multi-layer perceptrons (Liu et al., 2024b,a) or Q-formers (Li et al., 2023). Finally, the LLM generates responses by combining the visual and textual inputs.

2.2 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) is a technique designed to enhance the capabilities of LLMs by integrating external knowledge sources (Gao et al., 2023b; Gupta et al., 2024). It addresses key challenges of LLMs, such as hallucinations (Huang et al., 2023; Tonmoy et al., 2024; Xu et al., 2024; Fan et al., 2024) and time-sensitive information (Mousavi et al., 2024), by incorporating relevant information retrieved from external databases

or documents. The RAG process typically involves three main steps: retrieval, generation, and augmentation. First, a retriever identifies and extracts relevant document chunks from a knowledge base based on semantic similarity to the user’s query. These retrieved chunks are then combined with the original query to form an augmented context, which is used as input for the LLM to generate a response.

2.3 LLM Citation Generation

LLM Citation Generation has gained attention as a way to enhance the verifiability and transparency of model-generated responses by including citations linked to external evidence. Citation generation improves the factual accuracy of LLMs answers (Gao et al., 2023a) and allows users to trace the sources of information, thereby increasing the credibility and explainability of outputs (Tahaei et al., 2024). Early work like ALCE (Gao et al., 2023a) proposed foundational methods and evaluation metrics for enabling LLMs to generate citations. Subsequent studies have improved citation quality through fine-tuning approaches (Huang et al., 2024; Li et al., 2024; Ye et al., 2024) or multi-stage pipelines (Zhang et al., 2024b; Hennigen et al., 2023; Lee et al., 2023). Although existing works have explored the citation generation for LLMs, no prior studies have addressed the problem of contextual image reference in multimodal settings or proposed corresponding solutions.

3 Problem Definition.

In this section, we present the problem definition for contextual image reference. Given an input prompt consisting of an ordered sequence of mixed images and text, denoted as $\{E_1, E_2, \dots, E_m\}$, where each element E_i is either an image I_i from the set $I = \{I_1, I_2, \dots, I_n\}$ or a text segment T_i from the set $T = \{T_1, T_2, \dots, T_k\}$, the model’s goal is to generate a textual response R that meets the prompt’s requirements. The response R can be with some images referenced through a contextual image ID that falls within the range $[1, n]$, corresponding to the images in the input set. These image references must be contextually appropriate and align with the textual descriptions or contextual requirements specified in the prompt, ensuring the response is coherent, relevant, and adheres to the prescribed format.

4 Method

To enhance the contextual image referencing capability of vision-language models, it is critical to incorporate contextually relevant image-text data during the model’s training phase. Specifically, datasets that feature contextually integrated image references should be included in the supervised fine-tuning (SFT) stages of the model’s development. We construct the training dataset by incorporating image references from the retrieved documents into appropriate positions within the original text-only responses generated by the LLM. To achieve this goal, two challenges need to be addressed. The first is to collect an SFT dataset containing responses with contextual image references. The second is to effectively fine-tune existing VLMs.

Overview. To address the initial challenge, we developed a method to generate responses with contextual image references using existing LLMs, VLMs and user prompts. Given a prompt, we first generate a text-based response, then create a caption for the input image using the VLM, and integrate this caption into the text. Our model training follows the VLM’s standard SFT approach, but we enhance image understanding by requiring the model to refer explicitly to input images. To preserve the general VLM capabilities, we combine the original SFT data with interleaved multi-image SFT data. The ImageRef-VL framework is illustrated in Figure 2.

4.1 Training Data Construction

To construct the training dataset, we use existing LLMs and VLMs to create a high-quality dataset through a multi-stage few-shot learning approach, and then manually select qualified samples for model training. This method significantly reduces the labeling effort. Specifically, the data generation process involves three steps: generating text-based responses from pure text content, generating captions for each image based on text context, and adding the images references into the responses.

Text Response Generation. In this step, we remove the images from the prompt and provide some reference text for the model to generate a pure text response.

Image Caption Generation. In this step, we need to generate context-based image descriptions. The image description should not only include information about the image itself but also complement the

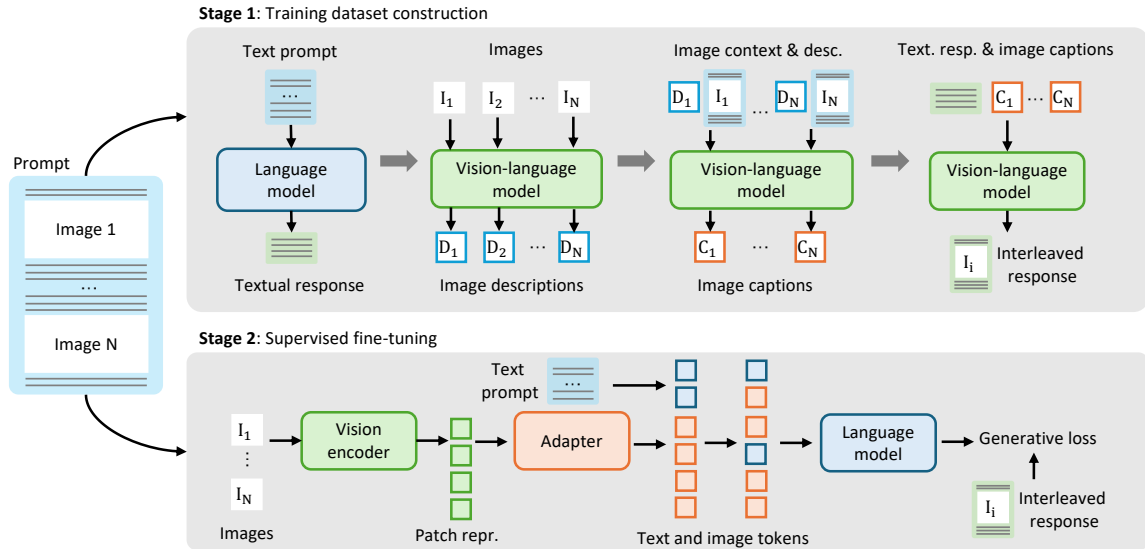


Figure 2: **The training strategy of the proposed IMI-VL model.** Stage 1: Training dataset construction involves generating textual responses and image descriptions through a language model and a vision-language model. These are combined into interleaved responses using image contexts and captions. Stage 2: Supervised fine-tuning refines the model with a vision encoder, adapter, and language model, optimizing through generative loss.

information related to the image described in the text. For example, when describing an image in the Wikipedia page about Einstein, it is not enough to simply say, ‘This is a portrait of an elderly man with white hair’; we also need to complete the information by adding, ‘This is a portrait of Einstein in his later years.’ However, directly using the existing VLM to perform this task can lead to over-reliance on context: when the image is unrelated to the context, the model may incorrectly force a connection, and some contextual information may be left incomplete.

To address this issue, we propose a two-stage image description generation approach. In the first stage, we generate a description based solely on the image itself. In the second stage, we provide the image description along with the context to the VLM, asking it to supplement the missing information. We use in-context learning and include examples in both stages to guide the VLM on how to better generate the image description and complete the information.

Image Insertions. We input the generated image captions and text responses into an LLM, asking the model to insert the images into the text response. The final results, after manual filtering, form our IMI-interleave training dataset.

Mixture of Datasets. To ensure the LLM truly understands the images in context, we also incorporate a certain proportion of contextual image

caption generation tasks in the training set, where VLMs generate captions for multiple images in the prompt based on the context, with the captions being those produced in the second step. Additionally, to prevent a significant decline in the performance of VLMs on other image-related tasks, we mix in a proportion of the InternVL2 SFT dataset into our training set.

4.2 Model Training

We used the constructed dataset to perform supervised fine-tuning on the VLM to enhance its ability to understand interleaved multi-image tasks. In this subsection, we will briefly introduce the model architecture and training loss.

Model Architecture. As shown in stage 2 of Figure 2, our fine-tuned model consists of three components: the vision encoder, the adapter, and the language model. Given a user prompt and retrieved documents with images, the vision encoder processes each image, extracting patch features from the image. The adapter, acting as a bridge between the vision encoder and the language model, maps the patch features into the embedding space of the language model, resulting in image tokens. The textual part of the prompt is tokenized and embedded to obtain text tokens. Finally, the text tokens and image tokens are fed into the language model in their original positions for modeling.

Training Loss. We proceed the typical supervised

fine-tuning process (Ouyang et al., 2022). The dataset \mathcal{D} is composed of N prompt-response pairs: $\mathcal{D} = \{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^N$, where both prompt \mathbf{x}^i and ground-truth response \mathbf{y}^i are a sequence of tokens. We denote $p(\mathbf{y}_j^i | \mathbf{x}^i \oplus \mathbf{y}_{<j}^i)$ as the probability of the outputting next token as \mathbf{y}_j^i given previous tokens $\mathbf{x}^i \oplus \mathbf{y}_{<j}^i$, where \oplus is the concatenation operator and $\mathbf{y}_{<j}^i$ denotes the tokens before index j . The training loss is then $\mathcal{L} = -\log \prod_{j=1}^{n^i} p(\mathbf{y}_j^i | \mathbf{x}^i \oplus \mathbf{y}_{<j}^i)$, with n^i being the length of \mathbf{y}^i and the optimization variable being a VLM.

5 Evaluation Settings.

To evaluate the performance of model supporting contextual image references, we propose three kinds of evaluation metrics, i.e., textual content evaluation, image position evaluation, overall response evaluation.

5.1 Automated Evaluation

Text Evaluation. When a model generates responses with contextual image references, we can evaluate the textual portion of the response. Following the existing LLM-as-judge approach (Zheng et al., 2023), allowing a large language model to score the response and then calculating the candidate model’s average score across all test samples. Considering that current large language models may not perform well in understanding multi-image prompts, we only include the textual prompt and provide a reference answer for scoring.

Image Position Evaluation. Inspired by existing work on image position prediction (Muraoka et al., 2020), when a model generates contextual image references, we can evaluate whether each image is placed in an appropriate position. However, previous approaches rely on an existing image-inserted dataset, checking whether the model places the image in a specific, pre-defined position or if the image ranks within the top-K choices. However, this metric does not align well with the actual user experience, as multiple images within the candidate pool could be suitable for position i .

To address this issue, we redesigned the image position evaluation metric. Specifically, for each potential image insertion point, we classify all images into four categories:

- 3-point images: Images that perfectly match the current contextual content.

- 2-point images: Images that match the current contextual content but are of low quality (e.g., blurry, obstructed).
- 1-point images: Images related to the main subject mentioned in the current context.
- 0-point images: Images not related to the current contextual content.

Finally, we can obtain a testing dataset in the following format:

$$\mathcal{D}_{\text{test}} = \left\{ p_i : \left\{ s : [I_{sp_i}^1, \dots, I_{sp_i}^{M_{sp_i}}] \mid s \in [0, 3] \right\} \mid p_i \in P_i, i \in [1, N] \right\}, \quad (1)$$

where P_i is all potential image insertion points of the i -th testing sample, $I_{sp_i}^j$ is the j -th image with s point for position p_i , M_{sp_i} is the maximum number of s -point images at position p_i , and N is the number of samples in the testing dataset.

Based on this label definition, we designed three metrics: Precision, Recall3, F1 and Score.

- **Precision** is defined as the accuracy of the illustrations, meaning the probability that an inserted image is a nonzero score image. It is defined as follows:

$$\text{Precision} = \frac{\sum_{i=1}^N \mathbb{I}(s_i > 0)}{N}, \quad (2)$$

where N is the total number of inserted images, and $\mathbb{I}(s_i > 0)$ is an indicator function that returns 1 if the score s_i of the i -th image is greater than zero, and 0 otherwise.

- **Recall3** is defined as the coverage rate of the relevant images, representing the probability that the images were inserted at all possible positions where the 3-point images could be inserted. Since the same image cannot be inserted in two different positions, we used the BPM algorithm to calculate the maximum number of relevant images that can be inserted under the current sample’s score label. The Recall metric is defined as follows:

$$\text{Recall3} = \frac{\sum_{p \in P} \mathbb{I}(s_p = 3)}{\sum_{p \in P} M_{3p}}, \quad (3)$$

where P is the set of all possible image insertion points, and $\mathbb{I}(s_p = 1)$ is an indicator function that returns 1 if an image with score 3 is inserted at position p , and 0 otherwise.

Table 1: Statistic details of datasets used in our experiments.

Train			
Dataset	# Sample	# Image	Avg. prompt len.
CIR-Interleave	7,645	73,833	5,481.67
CIR-Caption	5,633	29,558	2,215.84
InternVL2-SFT	1,267,819	1,227,131	501.54
Test			
Dataset	# Sample	# Image	Avg. prompt len.
CIR-Test	456	3,767	5,884.99

- **F1** is the harmonic mean of precision and recall3, which is formulated as follows:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall3}}{\text{Precision} + \text{Recall3}} \quad (4)$$

5.2 Human Evaluation

In addition to evaluating the text generation content and image placement, it is also necessary to assess the overall quality of the generated interleaved response. However, there is currently no effective automated method to evaluate the quality of a text-image interleaved response. Therefore, we employ human evaluation, assigning separate scores for the text, image, and overall response quality. For scoring in all three aspects, we use a 5-point Likert scale (Joshi et al., 2015) as the rating option.

6 Experimental Results

6.1 Experimental Settings

According to the description in Section 4.1, we constructed two datasets, the CIR-Interleave and CIR-Caption datasets, to enhance the model’s ability of contextual image referencing. The CIR-Interleave dataset consists of multiple user prompt and retrieved documents, and the model is required to generate responses with contextual image references based on these texts and user prompts. The CIR-Caption dataset contains contextual caption texts, and the model is tasked with providing contextually relevant captions for each image mentioned within the text. Additionally, we blend the InternVL SFT dataset at a certain ratio to fine-tune the model. During the testing phase, we construct the CIR-Test dataset as outlined in Section 5, and report evaluation scores on this dataset, including text evaluation scores, image position evaluation metrics such as precision, recall, F1, and human evaluation scores. Detailed statistics for all the datasets used during training and testing can be found in Table 1.

We conduct experiments on InternVL2-8B and InternVL2-26B (Chen et al., 2024), using 16 A100 GPUs, optimizing the full-fine-tuning training process with DeepSpeed Zero-3 and gradient checkpointing for improved memory efficiency and scalability. Both models are trained with a global batch size of 128, utilizing a learning rate of 4e-5 for the 8B model and 2e-5 for the 26B model, with corresponding weight decay values of 0.01 and 0.05. Training is conducted over 1000 steps, with a maximum sequence length of 16,384 tokens. For ImageRef-VL-26B, we mix the InternVL2-SFT dataset with the combined MI-Interleave and IMI-Caption datasets at a 1:1 ratio and train the model for 550 steps. For ImageRef-VL-8B, we mix the InternVL2-SFT dataset with the combined MI-Interleave and IMI-Caption datasets at a 1:4 ratio and train the model for 950 steps.

6.2 Performance Comparison

In the subsection, we compare the performance of ImageRef-VL with the baseline methods, including existing close-sourced vision-language models:

- **GPT-4o** (Hurst et al., 2024): a fast, cost-effective, multimodal large language model by OpenAI.
- **Claude-3.5-Sonnet** (Anthropic, 2024): A multimodal large language model by Anthropic, offering advanced safety and language capabilities.

and open-sourced vision-language models:

- **Phi-3.5-Vision** (Abdin et al., 2024): a lightweight, open multimodal large language model by Microsoft.
- **Qwen2-VL-7B** (Wang et al., 2024): the latest version of the vision language models in the Qwen model families.
- **InternVL2-26B** (Chen et al., 2024): the latest addition to the InternVL series of multimodal large language models.

and the three-stage response generation method introduced in Section 4.1 with GPT-4o and InternVL2-26B. We report the text evaluation score and human evaluation score for all models. For image position evaluation score, controlled sampling⁶ is used for open-sourced VLMs to complete image

⁶<https://github.com/dottxt-ai/outlines>

Table 2: Performance comparison of our ImageRef-VL with baselines. The top method for each metric is in **bold**, and the second-best is underlined.

Method	Model	Text eval.	Image position evaluation			Human evaluation		
		Score	Precision	Recall3	F1	Text	Image	Overall
Close-sourced VLMs	GPT-4o	7.27	—	—	—	2.61	3.78	3.17
	Claude-3.5-sonnet	<u>7.35</u>	—	—	—	2.60	3.91	3.32
Open-sourced VLMs	Phi-3.5-Vision	2.51	100.00	5.73	10.83	1.55	2.20	1.57
	Qwen2-VL-7B	1.95	100.00	5.73	10.83	1.82	2.15	1.73
	InternVL2-26B	4.84	79.59	10.56	18.65	1.92	2.51	1.87
Three-stage generation	GPT-4o	7.86	66.09	62.94	64.48	2.77	3.80	3.25
	InternVL2-26B	5.85	59.05	<u>59.43</u>	<u>59.24</u>	2.09	3.21	2.29
Our method	ImageRef-VL-8B	7.09	65.20	37.25	47.41	2.97	<u>4.05</u>	3.52
	ImageRef-VL-26B	7.30	63.75	29.26	40.11	<u>2.90</u>	4.08	<u>3.36</u>

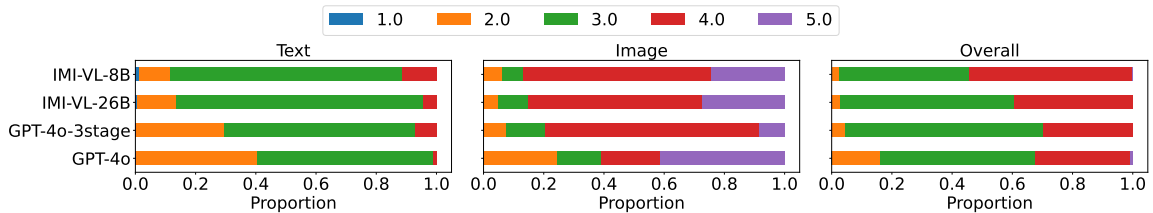


Figure 3: Human evaluation score distribution of four methods.

references based on the given text. For three-stage methods, we provide text in prompts and ask the model to insert image references (see Section 4.1 Image Insertions part). Since closed-source VLMs lack controlled sampling capabilities, their image position scores cannot be reported. The experimental results are shown in Table 2. We also present the detailed distribution of human evaluation scores for GPT-4o, GPT-4o three-stage, as well as our ImageRef-VL-8B and -32B in Figure 3.

Effectiveness of ImageRef-VL. In terms of the human evaluation metrics, ImageRef-VL-8B and ImageRef-VL-26B achieved the best performance in text scores, illustration scores, and overall experience. An 88% performance improvement is achieved over state-of-the-art open-source VLMs. This demonstrates the effectiveness of our approach. By further examining Figure 3, we can observe that the proportion of severe bad cases produced by ImageRef-VL is significantly lower than that of other methods, with the bad case rate of ImageRef-VL-26B being lower than that of ImageRef-VL-8B.

Open-source VLMs are significantly inferior to that of closed-source VLMs. Besides, among open-source VLMs, InternVL2 outperforms the others. Currently, open-source VLMs have not considered multi-image contextual understanding tasks

during the SFT phase, which could be a reason for their weaker performance. Additionally, there is an inherent performance gap between open-source VLMs and state-of-the-art closed-source VLMs. Among the three open-source VLMs we tested, InternVL2 considered contextual multi-image input during the pretraining phase, which might explain why it performs better.

Three-stage approach yields better results compared to direct end-to-end inference. The LLMs ability to understand long-context multi-image scenarios is weaker than its ability to perform text-based tasks. Additionally, the caption generation in the three-stage approach benefits from in-context learning, which provides a solid foundation for accurate illustration in the final stage.

Effectiveness of Automatic Metrics. In addition to manual evaluation metrics, we propose two types of automatic evaluation metrics to efficiently conduct preliminary evaluation and model screening. For text evaluation, the Pearson correlation with the text score of human evaluation is 0.9033 ($p < 0.005$), demonstrating its validity. For image position evaluation, due to differences in illustration processes between the three-stage and end-to-end approaches, a fair comparison is not feasible. Excluding the three-stage approach, we calculate the Pearson correlation between the F1 score

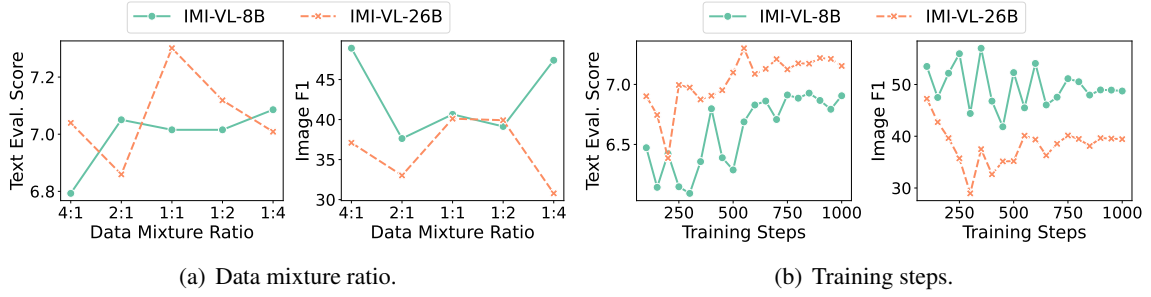


Figure 4: Impact of the hyper-parameters on our ImageRef-VL-8B and ImageRef-VL-26B.

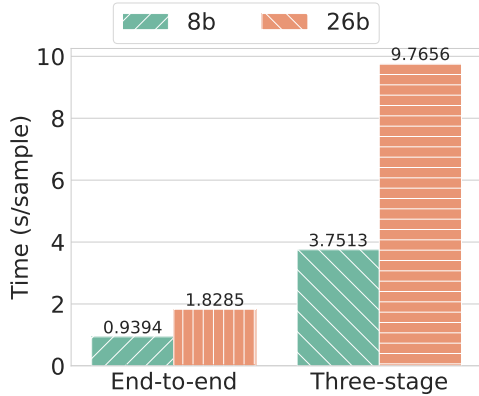


Figure 5: Computational cost comparison between our ImageRef-VL and three-stage methods.

and image score of human evaluation as 0.9854 ($p < 0.005$), confirming the metric’s validity.

6.3 Computation Overhead Analysis

In this section, we compare the computational costs of ImageRef-VL, an end-to-end interleaved response generation approach, with a three-stage generation approach. Define the number of images in a sample as N , the textual context length for each image as L , the total context length as M , the number of tokens occupied by a single image in the VLM as P , the response length as R , and the length of each image caption is C . The computational complexity of end-to-end method is

$$O((M + NP + R)^2), \quad (5)$$

while the complexity of three-stage method is

$$O((M + R)^2 + 9N(L + P + C)^2 + 9(R + NC)^2). \quad (6)$$

To more intuitively understand the difference in computational costs between the two methods, Figure 5 contrasts ImageRef-VL-8B and ImageRef-VL-26B with their respective three-stage

InternVL2-8B and InternVL2-26B counterparts. More specifically, we tested the end-to-end and 3-stage approaches of the 8B and 26B models on a machine equipped with 8 NVIDIA A100-SXM4-80GB GPUs to measure the execution time. The experimental results demonstrate that the end-to-end approach significantly reduces computational overhead compared to the three-stage scheme.

6.4 Hyper-parameter Analysis

Throughout the experiment, we focus on the impact of two hyper-parameters on the results: the proportion of mixed InternVL2 SFT data and the training steps of the model.

Data Mixture Ratio. Figure 4(a) shows the results for different data mixture ratios. For ImageRef-VL-26B, both the text evaluation score and image position evaluation F1 score initially rise and then decline, peaking at around a 1:1 data mixture ratio. In contrast, ImageRef-VL-8B’s text evaluation score steadily increases with a higher MI data proportion, while the image position evaluation F1 score first drops and then rises.

Training Steps. The results of different data mixture ratio are shown in Figure 4(b). For ImageRef-VL-26B, the text evaluation score of the model increases with the number of training steps, converging and stabilizing around 500 steps. The F1 score of image position evaluation initially fluctuates but also stabilizes around 500 steps. ImageRef-VL-8B shows a similar trend, but the convergence occurs around 700 steps.

7 Conclusion

In this paper, we introduced *Contextual Image Reference* as a novel capability for Vision-Language Models and presented ImageRef-VL, a framework that significantly advances the state-of-the-art in this domain. Through our carefully curated training data and proposed fine-tuning approach, we

demonstrated substantial improvements in VLMs’ ability to incorporate relevant images contextually in their responses. Our comprehensive evaluation framework, including both automated metrics and human assessment, validates the effectiveness of our approach. ImageRef-VL demonstrates superior performance over baseline models, achieving significantly better contextual image referencing while being computationally more efficient than multi-stage approaches. Our end-to-end system advances the development of visually-aware conversational AI. Future work could explore dynamic image generation and modification capabilities. We believe this work provides a strong foundation for research in multimodal AI systems with enhanced visual understanding.

Limitations

Our work validates the capability of enhanced VLMs to perform interleaved multi-image understanding, thereby exploring the feasibility of contextual image reference. However, our current experiments are based on post-finetuning of InternVL2. Starting from a well-pretrained VLM and incorporating the collected dataset into the supervised finetuning phase might yield better results. Additionally, the current model still exhibits a probability of bad cases. On one hand, collecting more and richer training datasets might address this issue; on the other hand, designing rewards and leveraging techniques like RLHF or RLAIIF could be employed for further fine-tuning of the model.

References

- Marah Abidin, Jyoti Aneja, Hany Awadallah, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Anthropic. 2024. [Claude 3.5 sonnet](#).
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *CVPR*, pages 24185–24198.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *KDD*, pages 6491–6501.
- Constanza Fierro, Reinald Kim Amplayo, Fantine Huot, Nicola De Cao, Joshua Maynez, Shashi Narayan, and Mirella Lapata. 2024. Learning to plan and generate text with citations. *arXiv preprint arXiv:2404.03381*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2024. [Mme: A comprehensive evaluation benchmark for multimodal large language models](#). *Preprint*, arXiv:2306.13394.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023a. Enabling large language models to generate text with citations. *arXiv preprint arXiv:2305.14627*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023b. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Shailja Gupta, Rajesh Ranjan, and Surya Narayan Singh. 2024. [A comprehensive survey of retrieval-augmented generation \(rag\): Evolution, current landscape and future directions](#). *Preprint*, arXiv:2410.12837.
- Lucas Torroba Hennigen, Shannon Shen, Anirudha Nrusimha, Bernhard Gapp, David Sontag, and Yoon Kim. 2023. Towards verifiable text generation with symbolic references. *arXiv preprint arXiv:2311.09188*.
- Chengyu Huang, Zeqiu Wu, Yushi Hu, and Wenya Wang. 2024. Training language models to generate text with citations via fine-grained rewards. *arXiv preprint arXiv:2402.04315*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *TOIS*.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, pages 6700–6709.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar Pal. 2015. Likert scale: Explored and explained. *British journal of applied science & technology*, 7(4):396–403.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *NIPS*, 35:22199–22213.

- Dongyub Lee, Taesun Whang, Chanhee Lee, and Heuseok Lim. 2023. Towards reliable and fluent large language models: Incorporating feedback learning loops in qa systems. *arXiv preprint arXiv:2309.06384*.
- Dongfang Li, Zetian Sun, Baotian Hu, Zhenyu Liu, Xinshuo Hu, Xuebo Liu, and Min Zhang. 2024. Improving attributed text generation of large language models via preference learning. *arXiv preprint arXiv:2403.18381*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742. PMLR.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *CVPR*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024b. Visual instruction tuning. *NIPS*, 36.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. 2024. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.
- Meta. 2024. [Llama 3.2](#).
- Seyed Mahed Mousavi, Simone Alghisi, and Giuseppe Riccardi. 2024. Is your llm outdated? benchmarking llms & alignment algorithms for time-sensitive knowledge. *arXiv preprint arXiv:2404.08700*.
- Masayasu Muraoka, Ryosuke Kohita, and Etsuko Ishii. 2020. Image position prediction in multimodal documents. In *LREC*, pages 4265–4274.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *NIPS*, 35:27730–27744.
- Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. 2022. Synchronesh: Reliable code generation from pre-trained language models. *arXiv preprint arXiv:2201.11227*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763.
- Jiajun Shen, Tong Zhou, Suifeng Zhao, Yubo Chen, and Kang Liu. 2024. Citekit: A modular toolkit for large language model citation generation. *arXiv preprint arXiv:2408.04662*.
- Marzieh Tahaei, Aref Jafari, Ahmad Rashid, David Alfonso-Hermelo, Khalil Bibi, Yimeng Wu, Ali Ghodsi, Boxing Chen, and Mehdi Rezagholizadeh. 2024. Efficient citer: Tuning large language models for enhanced answer quality and verification. In *Findings of NAACL*, pages 4443–4450.
- SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313*.
- A Vaswani. 2017. Attention is all you need. *NIPS*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817*.
- Xi Ye, Ruoxi Sun, Sercan Arik, and Tomas Pfister. 2024. Effective large language model adaptation for improved grounding and citation generation. In *NAACL*, pages 6237–6251.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *ICML*, pages 41092–41110.
- Jiajie Zhang, Yushi Bai, Xin Lv, Wanjuan Gu, Danqing Liu, Minhao Zou, Shulin Cao, Lei Hou, Yuxiao Dong, Ling Feng, et al. 2024a. Longcite: Enabling llms to generate fine-grained citations in long-context qa. *arXiv preprint arXiv:2409.02897*.
- Jingyu Zhang, Marc Marone, Tianjian Li, Benjamin Van Durme, and Daniel Khashabi. 2024b. Verifiable by design: Aligning language models to quote from pre-training data. *arXiv preprint arXiv:2404.03862*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *NIPS*, 36:46595–46623.