

Benchmarking Retrieval-Augmented Generation in Multi-Modal Contexts

Zhenghao Liu¹, Xingsheng Zhu¹, Tianshuo Zhou¹, Xinyi Zhang¹, Xiaoyuan Yi²,
Yukun Yan³, Yu Gu¹, Ge Yu¹, Maosong Sun^{3*}

¹Department of Computer Science and Technology, Northeastern University, China

²Microsoft Research Asia, Beijing, China

³Department of Computer Science and Technology, Institute for AI, Tsinghua University, China
Beijing National Research Center for Information Science and Technology, China

Abstract

This paper introduces **Multi-Modal Retrieval-Augmented Generation (M²RAG)**, a benchmark designed to evaluate the effectiveness of Multi-modal Large Language Models (MLLMs) in leveraging knowledge from multi-modal retrieval documents. The benchmark comprises four tasks: image captioning, multi-modal question answering, multi-modal fact verification, and image reranking. All tasks are set in an open-domain setting, requiring RAG models to retrieve query-relevant information from a multi-modal document collection and use it as input context for RAG modeling. To enhance the context utilization capabilities of MLLMs, we also introduce **Multi-Modal Retrieval-Augmented Instruction Tuning (MM-RAIT)**, an instruction tuning method that optimizes MLLMs within multi-modal contexts. Our experiments show that MM-RAIT improves the performance of RAG systems by enabling them to effectively learn from multi-modal contexts. All data and code are available at <https://github.com/NEUIR/M2RAG>.

1 Introduction

With the rapid development of Large Language Models (LLMs), such as GPT-4 (OpenAI, 2023) and LLaMA (Touvron et al., 2023), they have demonstrated strong emergent abilities in many NLP tasks (Wei et al., 2022; Zhao et al., 2023). However, LLMs often face the issue of hallucinations, making them produce unreliable responses (Ji et al., 2023; Huang et al., 2023; Shuster et al., 2021). Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Asai et al., 2024b; Shi et al., 2024; Yao et al., 2023) has proven effective in mitigating this hallucination problem by integrating external knowledge with LLMs.

To enhance LLMs with retrieved knowledge, existing approaches typically feed retrieved documents as input contexts, prompting LLMs to

* indicates corresponding author.

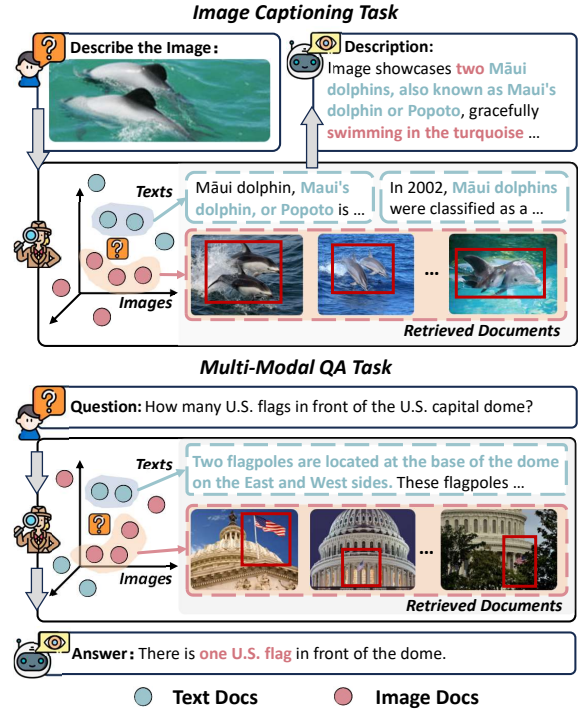


Figure 1: Illustration of Multi-Modal RAG Tasks. We incorporate multi-modal retrieval documents as the input context for MLLMs.

generate responses based on this in-context information (Ram et al., 2023). Existing RAG approaches (Petroni et al., 2021; Lin et al., 2024) usually focus on retrieving textual knowledge from corpora to aid LLMs in answering queries. Recent studies (Hu et al., 2024; Sharifmoghaddam et al., 2024) have extended RAG to Multi-modal Large Language Models (MLLMs), enabling them to address knowledge-intensive and information-seeking tasks that involve visual queries. However, these approaches largely rely on text or images as the sole sources of external knowledge, often overlooking the critical role of multi-modal data in providing richer and more comprehensive information that leads to producing more accurate answers (Hu et al., 2024; Liu et al., 2023b).

To advance RAG modeling in multi-modal scenarios, we introduce the Multi-Modal RAG (M²RAG) benchmark, designed to explore the effectiveness of MLLMs by feeding multi-modal retrieved documents as the input contexts to answer the question. As shown in Figure 1, we can use images or text as queries to retrieve multi-modal documents via multi-modal dense retrievers (Liu et al., 2023b; Zhou et al., 2024b,a). These multi-modal documents are then used as the input contexts to assist MLLMs during generation. Specifically, our M²RAG benchmark includes four distinct tasks: image captioning, multi-modal question answering, multi-modal fact verification, and image reranking. Different from existing works (Aghajanyan et al., 2022; Sharifymoghaddam et al., 2024), M²RAG is built upon high-quality datasets (Chang et al., 2022; Mishra et al., 2022) and designs four evaluation tasks in the open-domain setting, aiming to assess the effectiveness of MLLMs in leveraging the knowledge from multi-modal contexts.

In this paper, we also propose the **Multi-Modal Retrieval Augmented Instruction Tuning (MM-RAIT)** method to adapt MLLMs to the multi-modal in-context learning scenario, enhancing the effectiveness of MLLMs in utilizing the knowledge from these multi-modal retrieval documents. Specifically, we design task-specific prompt templates and fine-tune MLLMs on the M²RAG benchmark, making MLLMs maintain contextual awareness during generation. Our experimental results demonstrate that using retrieved knowledge significantly enhances MLLMs’ performance, achieving significant improvements in both zero-shot and few-shot settings. After training with MM-RAIT, MiniCPM-V and Qwen2-VL show an average improvement of 27% and 34% over vanilla RAG modeling methods, showing the effectiveness of MM-RAIT.

2 Related Work

Existing RAG models (Shi et al., 2024; Asai et al., 2024a; Yu et al., 2023b; Yan et al., 2024) typically rely on dense retrievers (Karpukhin et al., 2020; Xiong et al., 2021a; Ren et al., 2021; Xiong et al., 2021b; Gao and Callan, 2022) or sparse retrievers like BM25 (Robertson et al., 2009) for text document retrieval. More recent efforts have integrated multi-modal retrieval methods, allowing the inclusion of rich external knowledge from different modalities within RAG frameworks. For example, some works (Liu et al., 2023b; Zhou

et al., 2024b) have introduced unified multi-modal retrieval systems that map images and texts into a shared semantic space. This approach allows for single-modal matching, cross-modal matching, and modality routing within the embedding space. VISTA (Zhou et al., 2024a) further enhances multi-modal retrieval by optimizing synthetic training data and refining training strategies. These advancements enable the retrieval of multi-modal knowledge, providing a way for evaluating the effectiveness of MLLMs in multi-modal contexts.

Multi-modal Large Language Models (MLLMs) (Achiam et al., 2023; Team et al., 2023; Sun et al., 2024b,a; Aghajanyan et al., 2022; Lu et al., 2024) have proven their effectiveness in understanding, integrating, and utilizing both visual and textual knowledge in generation tasks. Models like BLIP (Li et al., 2022, 2023), LLaVA (Liu et al., 2023a), and Flamingo (Alayrac et al., 2022) build the MLLMs by combining pre-trained vision encoders with Large Language Models (LLMs) to process multi-modal inputs for generation. Thriving on the advancements in MLLMs, researchers pay more attention to extending the advantages of Retrieval-Augmented Generation (RAG) to these MLLMs, enhancing their generation capability.

Multi-modal RAG has demonstrated its potential to enhance knowledge-intensive and information-seeking tasks, such as question answering (Chang et al., 2022; Marino et al., 2019) and fact verification (Mishra et al., 2022). These models utilize retrieval-based multi-modal documents to provide richer and contextually relevant information. Additionally, other works have applied multi-modal RAG to improve the performance of MLLMs on the tasks like image captioning (Lin et al., 2014; Young et al., 2014) and generation (Yasunaga et al., 2023; Yu et al., 2023a; Sharifymoghaddam et al., 2024). However, existing multi-modal benchmarks (Johnson et al., 2017; Schuhmann et al., 2021; Lin et al., 2014; Young et al., 2014; Marino et al., 2019) are typically tailored to specific tasks and lack a comprehensive framework for evaluating multi-modal RAG systems.

3 M²RAG Benchmark for Multi-Modal Retrieval-Augmented Generation

In this section, we describe our Multi-Modal Retrieval-Augmented Generation (M²RAG) benchmark. We first introduce the RAG tasks in M²RAG,

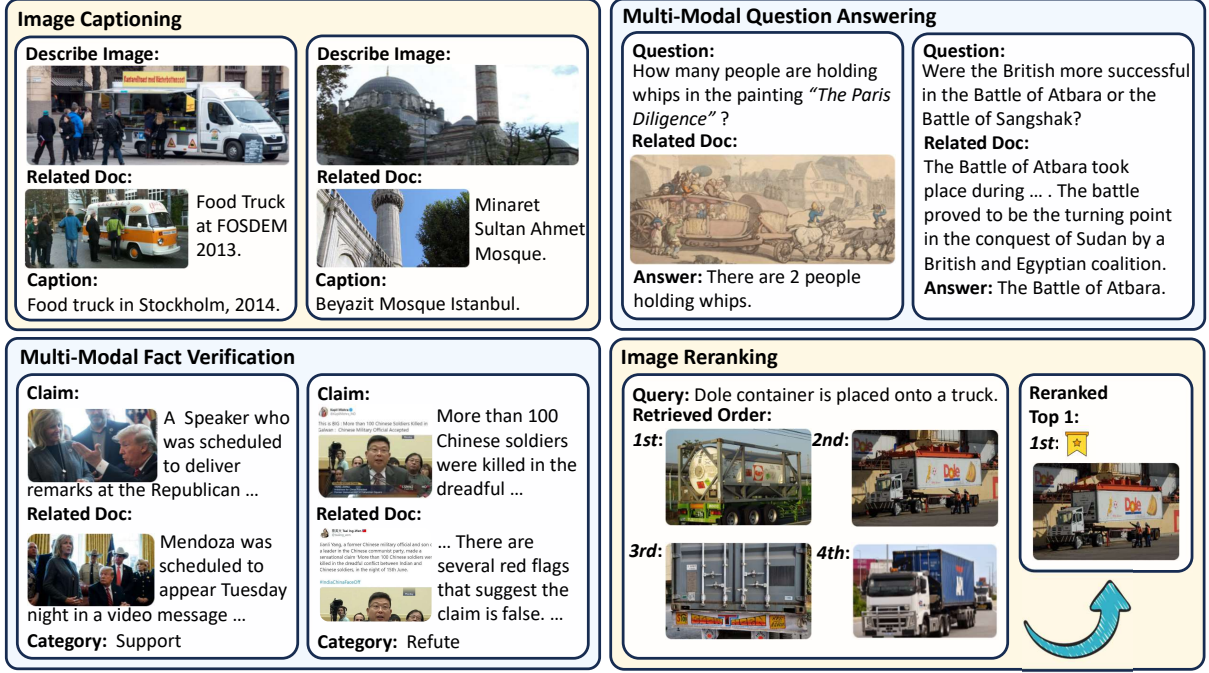


Figure 2: Examples of Different Tasks Defined in the M²RAG Benchmark. All tasks are designed for the open-domain setting. Thus, we present the input, ground truth answers, and retrieved documents for each task.

followed by a detailed explanation of the construction process. Finally, we present a comparison of existing multi-modal benchmarks with M²RAG.

Task Definition. As shown in Figure 2, M²RAG defines four tasks to evaluate the capabilities of MLLMs in open-domain RAG scenarios: image captioning, multi-modal question answering, multi-modal fact verification, and image reranking. For each task, MLLMs are required to retrieve knowledge from the multi-modal document collection \mathcal{D} and generate responses to answer the question q . Details of the prompt templates of different tasks are shown in Appendix A.3.

Image Captioning Task. Image Captioning is a widely used task for evaluating the performance of multi-modal RAG models (Aghajanyan et al., 2022; Sharifymoghaddam et al., 2024). In this task, an image is provided as the query q , and the document collection \mathcal{D} is constructed using image documents that contain captions. The goal of image captioning is to generate concise and semantically coherent captions that accurately describe the image content. Unlike previous works (Aghajanyan et al., 2022; Sharifymoghaddam et al., 2024), we source image captions from WebQA (Chang et al., 2022), where all image documents are collected from Wikimedia Commons. These captions often include important details, such as named entities, which make the

task more challenging and provide crucial information for query matching (Liu et al., 2023b). More comparison details of the image captioning tasks in different benchmarks are shown in Appendix A.4.

Multi-Modal Question Answering Task. Following the WebQA benchmark (Chang et al., 2022), the Multi-Modal QA task involves answering text-based queries q by leveraging both text and image documents. The document collection \mathcal{D} consists of both text and image documents containing captions. Additionally, we extend WebQA to an open-domain setting by instructing the retriever to return query-relevant documents from the collection \mathcal{D} , as demonstrated by Liu et al. (2023b).

Multi-Modal Fact Verification Task. The Multi-Modal Fact Verification task challenges MLLMs to verify the accuracy of claims using multi-modal evidence. In this task, the query q can be a multi-modal claim, and the document collection \mathcal{D} consists of both text and image documents, where the image documents do not contain captions. The relationship between the claim and the evidence is categorized into three possible outcomes: “Support”, “Refute”, or “Insufficient”, indicating whether the evidence supports, refutes, or lacks sufficient information to verify the claim. We build this task on the Factify dataset (Mishra et al., 2022), but we focus on open-domain fact verification by re-

Benchmarks	Input Modality	Knowledge Source	Retrieval Modality	Multi Task	Open Domain
MSCOCO (2014)	Image	Web	✗	✗	✗
Flickr30K (2014)	Image	Web	✗	✗	✗
K-VQA (2019)	Multi	Wikipedia	✗	✗	✗
WebQA (2022)	Text	Wikipedia	Multi	✗	✗
MRAG-Bench (2024)	Multi	Web and Other Sources	Image	✗	✗
M²RAG (Ours)	Multi	Wikipedia, Twitter	Multi	✓	✓

Table 1: Comparison of Multi-Modal Benchmarks.

trieving evidence from a multi-modal document collection (Thorne et al., 2018).

Image Reranking Task. In the Image Reranking task, the objective is to identify the most relevant images based on a given image description. Here, the image description serves as the query q , and the document collection \mathcal{D} consists of image documents that do not include captions. For each description, we first use a multi-modal retriever to retrieve image candidates based solely on their image features and then rerank the images using MLLMs. To adapt MLLMs for this task, we follow prior work (Muennighoff, 2022) and compute the Perplexity (PPL) score for reranking image candidates based on their image features. This approach models the relevance between queries and images in a manner similar to image captioning, where a lower PPL score indicates greater relevance between the candidate image and the given query.

Details of Data Construction. To build the M²RAG benchmark, we collect data from two widely used datasets, WebQA (Chang et al., 2022) and Factify (Mishra et al., 2022), constructing 3,000 instances for both training and evaluation processes for each task.

First, we select WebQA to build the tasks for image captioning, multi-modal QA, and image reranking. For the multi-modal QA task, we select an equal number of text-based and image-based QA pairs from WebQA to construct the dataset. For image captioning and image reranking tasks, we randomly select image-text pairs with a similarity score greater than 0.65, which are then split into training and evaluation sets. The retrieval corpus for these tasks consists of all image-text pairs except the selected ones. Then, for the multi-modal fact verification task, we use the Factify dataset and the entire set of image-text documents in the dataset is used as the retrieval corpus. Additionally, we consolidate the modality based categories into three: “Support”, “Refute”, and “Insufficient.”

Benchmark Comparison. The comparison of different benchmarks is presented in Table 1.

Existing multi-modal benchmarks primarily evaluate the effectiveness of MLLMs on single tasks, such as image captioning (Lin et al., 2014; Young et al., 2014) or question answering (Chang et al., 2022), limiting their ability to conduct comprehensive evaluations of MLLMs in multi-modal RAG scenarios. In contrast, M²RAG offers several unique features for a more thorough evaluation: 1) M²RAG defines four tasks that assess an MLLM’s ability to effectively understand and utilize retrieved knowledge. These tasks require MLLMs to perform reasoning and information matching based on both queries and contextual knowledge. 2) M²RAG incorporates the retrieval results as the multi-modal contexts for model inputs, avoiding the need for separate processing of the retrieval documents of different modalities. 3) M²RAG adapts these tasks to an open-domain setting, requiring MLLMs to retrieve knowledge from a comprehensive multi-modal document collection, offering a more realistic RAG scenario.

4 Instruction Tuning for Multi-Modal Retrieval-Augmented Generation

In this section, we present our Multi-Modal Retrieval-Augmented Instruction Tuning (MM-RAIT) method. First, we describe the framework for multi-modal Retrieval-Augmented Generation (RAG) (Sec. 4.1). Then, we introduce multi-task instruction tuning to enhance the performance of MLLMs in multi-modal RAG tasks (Sec. 4.2).

4.1 The Framework of Multi-Modal Retrieval-Augmented Generation

Given a query q , multi-modal RAG models first employ a retriever to search for query-relevant multi-modal documents \mathcal{D} and then feed these documents to MLLMs to assist them in answering the query q . Each document $d \in \mathcal{D}$ can be either an image document or a text document. The multi-modal RAG framework consists of two main components: the multi-modal retrieval module and the retrieval-augmented generation module.

Multi-Modal Retrieval. To retrieve documents from the multi-modal document collection \mathcal{D} , existing methods typically rely on multi-modal dense retrieval models (Zhou et al., 2024b,a).

Given a query q and a multi-modal document d , multi-modal dense retrieval models, such as VISTA (Zhou et al., 2024a), encode both as representations h_q and h_d , respectively, and map them into an embedding space for retrieval:

$$h_q = \text{Enc}(q); h_d = \text{Enc}(d), \quad (1)$$

where Enc denotes the encoder model. The query q can be either text or an image, and the multi-modal document d can be a text document or an image document. For documents containing captions, both image features and image captions are fed into the encoder model.

Next, we compute the similarity score $S(q, d)$ between the representations h_q and h_d of the query and document:

$$S(q, d) = \text{Sim}(h_q, h_d), \quad (2)$$

where Sim denotes cosine similarity. We then perform a KNN search (Johnson et al., 2019) to retrieve the top- k most relevant multi-modal documents $\tilde{\mathcal{D}} = \{d_1, \dots, d_k\}$ to the query q . During retrieval, the multi-modal retriever needs to conduct single-modality matching, cross-modality matching and modality routing in the embedding space (Liu et al., 2023b).

Multi-Modal RAG Module. After retrieval, we input the retrieved documents $\tilde{\mathcal{D}}$ and query q into the MLLM (\mathcal{M}), such as MiniCPM-V (Yao et al., 2024) or Qwen2-VL (Wang et al., 2024), to generate the output y :

$$y = \mathcal{M}(\tilde{\mathcal{D}}, q). \quad (3)$$

These retrieved documents provide external knowledge, which helps to update the parametric memory of the MLLM, enabling it to generate more accurate responses to the query q .

4.2 MM-RAIT: Multi-Task Multi-Modal Instruction Tuning for MLLMs

To adapt MLLMs to the multi-modal RAG scenario, we propose the **Multi-Modal Retrieval-Augmented Instruction Tuning** (MM-RAIT) method, designed to further enhance the performance of MLLMs across various RAG tasks.

To improve the MLLM generation process, we incorporate external knowledge to assist in answering the query (Eq.3). Specifically, we follow previous work (Ram et al., 2023) and concatenate the

representations of the retrieved documents $\tilde{\mathcal{D}}$ along with the query q as the input for the MLLM (\mathcal{M}) to generate the output y :

$$y = \mathcal{M}(\text{Instruct}_p, X(\tilde{\mathcal{D}}), q), \quad (4)$$

where Instruct_p is the instruction for the task p , and $X(\tilde{\mathcal{D}})$ denotes the concatenation of the representations of the retrieved documents:

$$X(\tilde{\mathcal{D}}) = X(d_1) \oplus \dots \oplus X(d_k). \quad (5)$$

For the i -th retrieved document d_i , its representation can be the text sequence for a text document, the image features for an image document, or the concatenation of both image features and caption for an image document that contains a caption.

Next, we gather queries from three tasks to form the query set Q : image captioning, multi-modal question answering, and multi-modal fact verification. For each query q in these tasks, the training objective for the model is to minimize the negative log-likelihood of generating the target sequence y^* :

$$\mathcal{L} = - \sum_{q \in Q} \sum_{t=1}^T \log P(y_t^* | y_{<t}^*, \tilde{\mathcal{D}}, q; \theta), \quad (6)$$

where T is the length of the ground truth response, y_t^* is the t -th token of the ground truth response, and θ represents the parameters of MLLM (\mathcal{M}).

5 Experimental Methodology

This section outlines the datasets, evaluation metrics, baselines, and implementation details used in our experiments.

Dataset. We use the M²RAG dataset to evaluate the performance of different MLLMs in the multi-modal RAG scenario. The dataset consists of four tasks: image captioning, multi-modal question answering, multi-modal fact verification, and image reranking. For multi-modal retrieval, we adopt VISTA (Zhou et al., 2024a), a universal embedding model to search for query-related documents. VISTA integrates image token embeddings into the BGE Text Embedding (Xiao et al., 2024) framework, enabling flexible processing of the inputs of both text and image data.

Evaluation Metrics. For image captioning and multi-modal QA tasks, we use BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and CIDEr (Vedantam et al., 2015) scores to assess performance. In the multi-modal fact verification task, we evaluate the performance of different RAG

Model	Image Captioning			Multi-Modal QA			MM Fact Verification		Image Reranking
	BLEU-4	ROUGE-L	CIDEr	BLEU-4	ROUGE-L	CIDEr	ACC	F1	FID↓
MiniCPM-V 2.6 (8B)									
Vanilla RAG	1.91	17.58	18.39	13.84	32.78	82.65	43.03	41.13	-
w/ top1	3.82	24.28	43.76	17.18	37.89	119.92	54.83	53.49	12.17
w/ top3	3.46	23.13	37.89	17.56	38.46	124.16	56.33	54.01	10.77
w/ top5	3.29	22.82	36.09	17.15	37.99	114.21	55.33	52.69	11.15
MM-RAIT	6.25	32.77	77.08	26.37	53.21	266.47	60.17	60.22	10.32
Qwen2-VL (7B)									
Vanilla RAG	2.24	19.48	26.01	18.77	39.05	153.40	45.43	34.23	-
w/ top1	3.79	25.43	46.32	21.21	42.14	187.99	51.60	41.05	12.17
w/ top3	3.62	25.70	45.08	20.98	41.90	178.28	52.43	41.94	9.88
w/ top5	3.62	25.31	44.45	21.26	42.41	181.56	52.00	41.64	9.71
MM-RAIT	10.53	39.79	123.97	32.00	62.49	329.05	65.13	62.97	9.15

Table 2: Overall Performance. We evaluate the performance of different RAG models implemented with MiniCPM-V 2.6 and Qwen2-VL on our M²RAG benchmark. MM-RAIT uses the top-5 multi-modal documents for inference.

models using accuracy (ACC) and F1 score. For the image reranking task, we use the Fréchet Inception Distance (FID↓) (Heusel et al., 2017)¹.

Baselines. We compare our models with two multi-modal baselines: MiniCPM-V 2.6 (Yao et al., 2024) and Qwen2-VL (Wang et al., 2024). These models are evaluated in a zero-shot setting to assess their effectiveness in leveraging multi-modal knowledge from the input context. We feed the top-1, top-3, and top-5 ranked documents into these MLLMs to evaluate their RAG performance.

Implementation Details. We apply Low-Rank Adaptation (LoRA) (Hu et al., 2022) to fine-tune both MiniCPM-V 2.6 and Qwen2-VL using the top-5 retrieved multi-modal documents for 2 epochs. The batch size is 4, with a maximum token limit of 4,096. A cosine learning rate scheduler is used, with the learning rate set to $1e-6$ for MiniCPM-V and $1e-4$ for Qwen2-VL. We fine-tune Qwen2-VL using LLaMA-Factory (Zheng et al., 2024) and set `max_pixels=512 × 512` for training and inference.

6 Evaluation Result

In this section, we first evaluate the performance of MLLMs on the M²RAG benchmark. We then conduct ablation studies to assess the impact of varying numbers of retrieved documents of different modalities. Following that, we analyze the role of different retrieval modalities in RAG models. Finally, case studies are shown.

6.1 Overall Performance

As shown in Table 2, we report the performance of various RAG models on the M²RAG benchmark. The vanilla RAG models directly use retrieved documents to augment LLMs, while MM-RAIT models fine-tune MLLMs within the RAG framework.

For these vanilla RAG models, performance generally improves as the number of retrieved documents increases. However, when retrieving the top-5 ranked documents, the overall performance of vanilla RAG models on most tasks is lower compared to using top-1 or top-3 documents. This suggests that vanilla LLMs still struggle to fully leverage multi-modal knowledge to enhance MLLMs. Although some related works also use image captioning tasks to evaluate RAG performance (Sharifymoghammad et al., 2024), the performance of these MLLMs on M²RAG is considerably worse, indicating that M²RAG offers a more challenging dataset for image captioning. In contrast to vanilla RAG models, both MiniCPM-V 2.6 and Qwen2-VL demonstrate strong performance across all tasks on the M²RAG benchmark after training with MM-RAIT. Specifically, MiniCPM-V 2.6 achieves an average improvement of over 27% across all tasks in M²RAG, while Qwen2-VL shows an even greater improvement of 34%. These results highlight the effectiveness of MM-RAIT, showcasing its ability to help MLLMs better utilize multi-modal contexts to enhance their performance.

6.2 Ablation Study

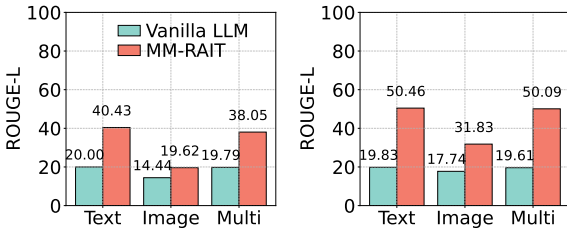
As shown in Table 3, we conduct ablation studies to evaluate RAG effectiveness with retrieved documents of different modalities and numbers. Specifically, we conduct two evaluation settings to evaluate the roles of different modalities: Only Text and Only Image. Only Text indicates removing all image features from multi-modal input contexts to enhance the MLLM, while Only Image removes all texts from top-ranked multi-modal documents.

Compared with the RAG models using top-3 ranked multi-modal documents for augmentation, the performance of vanilla RAG models usually de-

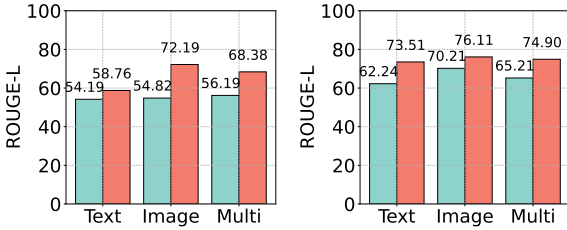
¹<https://github.com/mseitzer/pytorch-fid>

Model	#Doc	Image Captioning			Multi-Modal QA			MM Fact Verification		
		Only Text	Only Image	Multi	Only Text	Only Image	Multi	Only Text	Only Image	Multi
MiniCPM-V 2.6 (8B)										
MLLM	0	17.58	17.58	17.58	32.78	32.78	32.78	41.13	41.13	41.13
Vanilla RAG	3	24.20	17.15	23.13	38.25	37.57	38.46	53.21	43.61	54.01
	5	24.95	17.33	22.82	37.84	37.21	37.99	53.28	42.46	52.69
MM-RAIT	3	32.46	25.09	32.88	50.83	47.73	52.69	60.52	45.66	60.41
	5	33.04	25.38	32.77	50.32	47.12	53.21	60.88	47.71	60.22
Qwen2-VL (7B)										
MLLM	0	19.48	19.48	19.48	39.05	39.05	39.05	34.23	34.23	34.23
Vanilla RAG	3	26.87	19.12	25.70	41.02	43.10	41.90	43.16	33.56	41.94
	5	26.89	17.33	25.31	41.10	43.32	42.41	43.24	33.50	41.62
MM-RAIT	3	37.67	35.32	39.11	61.35	53.67	62.04	61.82	57.90	63.36
	5	38.39	35.27	39.79	61.95	53.97	62.49	61.40	57.41	62.97

Table 3: Ablation Study. We evaluate the performance of different retrieval modalities for candidate corpora on M²RAG benchmark. For Image Captioning and Multi-Modal QA, we use ROUGE-L as the evaluation metric. And F1-score is used for the MM Fact Verification task.



(a) MiniCPM Performance on Text Answerable Queries. (b) Qwen2 Performance on Text Answerable Queries.



(c) MiniCPM Performance on Image Answerable Queries. (d) Qwen2 Performance on Image Answerable Queries.

Figure 3: RAG Performance in Multi-Modal QA Task Using Retrieved Documents of Different Modalities. Text, Image, and Multi denote that retrieved text, image, and multi-modal documents are fed to different RAG models for evaluation.

creases with top-5 ranked documents, while MM-RAIT alleviates the performance decreases but also shows limited improvements. It illustrates that effectively using the multi-modal contextual knowledge is still challenging for existing MLLMs. Moreover, we further remove all texts or image features to show the roles of different modalities in RAG modeling. For all tasks, the RAG performance of the Only Text model slightly decreases, showing that these texts contribute to the primary knowledge source for these RAG models. After adding the image features, the RAG performance usually increases, showing that these image fea-

tures can improve the performance of RAG models. Even though different modalities show the effectiveness in multi-modal RAG modeling, it is still hard to effectively learn more crucial semantics from these image features to improve the RAG performance within multi-modal contexts.

6.3 Effectiveness of MLLMs in Different Modality-based RAG Scenarios

In this experiment, we investigate the impact of retrieved documents from different modalities on the effectiveness of RAG models.

As shown in Figure 3, we divide the multi-modal QA dataset of M²RAG into two groups: image-answerable queries and text-answerable queries. These categories represent queries that can be answered by image or text documents, respectively. We compare both vanilla RAG and MM-RAIT, implemented using MiniCPM-V and Qwen2-VL. Top-5 ranked documents from texts, images, and both modalities are fed to the different RAG models to evaluate the QA performance.

Figures 3(a) and 3(b) present the RAG performance on text-answerable queries. Overall, the RAG models using multi-modal retrieved documents exhibit comparable performance to those using only text-based documents, indicating that MLLMs can effectively learn from text documents to answer queries. Notably, vanilla RAG models show minimal differences in performance when using text, image, or both types of documents, whereas MM-RAIT significantly improves performance when leveraging documents from multiple modalities. This highlights the effectiveness of MM-RAIT in enabling MLLMs to learn from multi-modal contexts. Interestingly, vanilla MLLMs appear insensitive to the retrieved contexts, likely because they rely heavily on internal knowledge

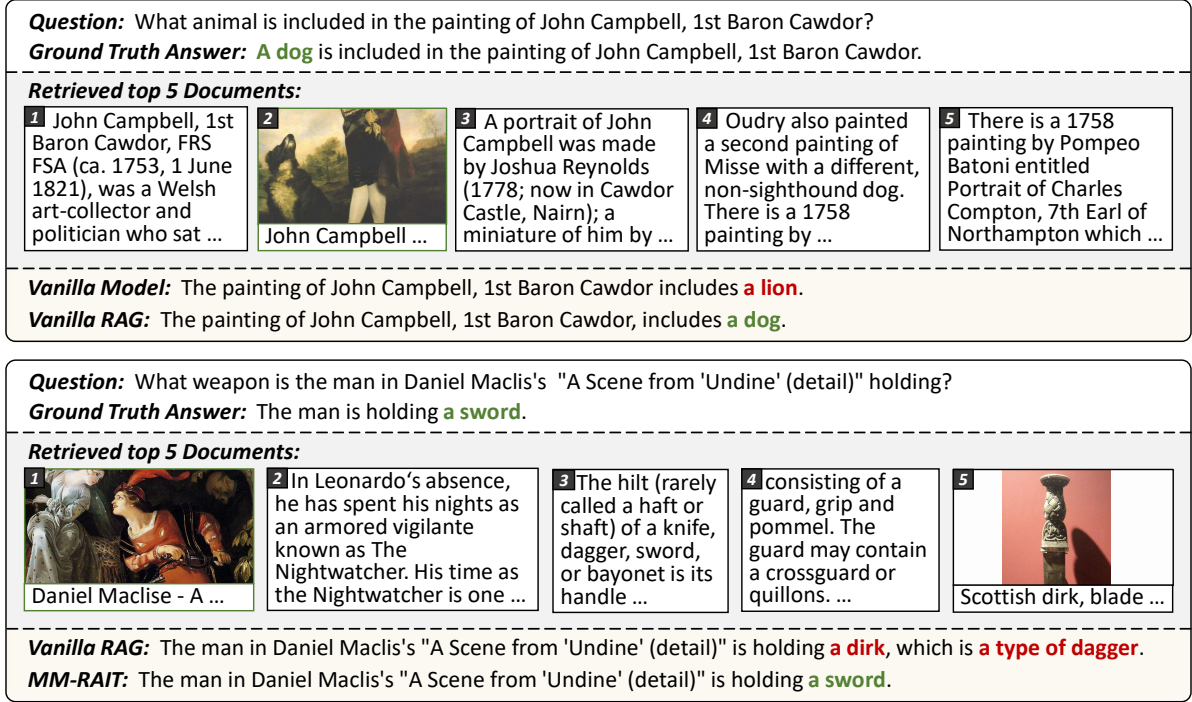


Figure 4: Case Studies. We highlight the **relevant phrases**, **correct answers**, and **query-unrelated phrases**.

when processing text-answerable queries.

Next, we evaluate the RAG performance on image-answerable queries, shown in Figures 3(c) and 3(d). The results indicate that RAG models using multi-modal documents generally outperform those using only text documents, confirming that incorporating image documents during retrieval enhances the ability of MLLMs to answer questions. The performance gap narrows for Qwen2-VL, suggesting that different MLLMs exhibit varying levels of reliance on multi-modal documents.

6.4 Case Study

In this section, we show two cases from Qwen2-VL in the Multi-Modal QA task of M²RAG to evaluate the effectiveness of the MM-RAIT method within the multi-modal retrieval contexts. More cases are shown in Appendix A.5.

As illustrated in Figure 4, in the first case, the question asks, “What animal is included in the painting of John Campbell, 1st Baron Cawdor?”. This requires the MLLM to match the “1st Baron Cawdor” and extract information about animals in the painting. Due to limited internal knowledge, the model encounters hallucination issues and generates an incorrect answer, “a lion”. When the retrieved multi-modal document of “1st Baron Cawdor” is fed into the MLLM, the vanilla RAG model can directly extract “dog” from the painting,

thus providing the correct response. This highlights the importance of multi-modal information in offering more intuitive and richer semantic insights to answer the question, underscoring the effectiveness of constructing the M²RAG benchmark.

In the second case, the question asks that, “What weapon is the man in Daniel Maclis’s A Scene from ‘Undine’ (detail) holding?” Based on retrieved documents, the vanilla RAG model focuses on the fifth document, which depicts a “Scottish dirk”. This leads the vanilla RAG model to generate an incorrect response, “holding a dirk”. After MM-RAIT training, the model can accurately identify the relevant document describing the man holding a sword and extract pertinent information from it, thereby generating the correct response.

7 Conclusion

This paper introduces **Multi-Modal Retrieval-Augmented Generation (M²RAG)**, a benchmark designed to evaluate MLLMs with retrieved multi-modal contexts across four tasks. To further enhance the utilization of retrieved information, we also propose a **Multi-Modal Retrieval-Augmented Instruction Tuning (MM-RAIT)** method, which optimizes MLLMs with multi-modal contexts as inputs, thereby improving their ability to effectively utilize retrieved information.

Limitations

Although our M²RAG benchmark includes four common multi-modal tasks, incorporating additional tasks can provide a more comprehensive evaluation of the capabilities of MLLMs. Furthermore, while MLLMs perform satisfactorily within retrieved multi-modal contexts, they still rely predominantly on textual data for some tasks. Finding ways to enable MLLMs to more effectively leverage multi-modal contexts remains a critical challenge that requires further exploration. Additionally, due to the performance limitations of multi-modal retrieval models, the quality of the retrieved multi-modal documents directly impacts the overall performance of MLLMs. Improving the accuracy of multi-modal retrieval remains a vital area for future research.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#). *ArXiv preprint*.
- Armen Aghajanyan, Bernie Huang, Candace Ross, Vladimir Karpukhin, Hu Xu, Naman Goyal, Dmytro Okhonko, Mandar Joshi, Gargi Ghosh, Mike Lewis, et al. 2022. [Cm3: A causal masked multimodal model of the internet](#). *ArXiv preprint*.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Proceedings of NeurIPS*.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024a. [Self-rag: Learning to retrieve, generate, and critique through self-reflection](#). In *Proceedings of ICLR*.
- Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi, and Wen-tau Yih. 2024b. [Reliable, adaptable, and attributable language models with retrieval](#). *ArXiv preprint*.
- Yingshan Chang, Guihong Cao, Mridu Narang, Jianfeng Gao, Hisami Suzuki, and Yonatan Bisk. 2022. [Webqa: Multihop and multimodal QA](#). In *Proceedings of CVPR*, pages 16474–16483.
- Luyu Gao and Jamie Callan. 2022. [Unsupervised corpus aware language model pre-training for dense passage retrieval](#). In *Proceedings of ACL*, pages 2843–2853.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. [Gans trained by a two time-scale update rule converge to a local nash equilibrium](#). In *Proceedings of NeurIPS*, pages 6626–6637.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *Proceedings of ICLR*.
- Wenbo Hu, Jia-Chen Gu, Zi-Yi Dou, Mohsen Fayyaz, Pan Lu, Kai-Wei Chang, and Nanyun Peng. 2024. [Mrag-bench: Vision-centric evaluation for retrieval-augmented multimodal models](#). *ArXiv preprint*.
- 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](#). *ArXiv preprint*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Computing Surveys*, (12):1–38.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. [Billion-scale similarity search with gpus](#). *IEEE Transactions on Big Data*, (3):535–547.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. 2017. [CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1988–1997.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of EMNLP*, pages 6769–6781.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Proceedings of NeurIPS*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. 2023. [BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models](#). In *Proceedings of ICML*, pages 19730–19742.

- Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. 2022. [BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation](#). In *Proceedings of ICML*, pages 12888–12900.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. [Microsoft coco: Common objects in context](#). In *Proceedings of ECCV*, pages 740–755. Springer.
- Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Richard James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. [RA-DIT: retrieval-augmented dual instruction tuning](#). In *Proceedings of ICLR*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. [Visual instruction tuning](#). In *Proceedings of NeurIPS*.
- Zhenghao Liu, Chenyan Xiong, Yuanhuiyi Lv, Zhiyuan Liu, and Ge Yu. 2023b. [Universal vision-language dense retrieval: Learning A unified representation space for multi-modal retrieval](#). In *Proceedings of ICLR*.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024. [Deepseek-vl: Towards real-world vision-language understanding](#).
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. [OK-VQA: A visual question answering benchmark requiring external knowledge](#). In *Proceedings of CVPR*, pages 3195–3204.
- Shreyash Mishra, S Suryavardan, Amrit Bhaskar, Parul Chopra, Aishwarya N Reganti, Parth Patwa, Amitava Das, Tanmoy Chakraborty, Amit P Sheth, Asif Ekbali, et al. 2022. [Factify: A multi-modal fact verification dataset](#). In *DE-FACTIFY@ AAAI*.
- Niklas Muennighoff. 2022. [Sgpt: Gpt sentence embeddings for semantic search](#). *ArXiv preprint*.
- R OpenAI. 2023. [Gpt-4 technical report](#). *arXiv*, pages 2303–08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of ACL*, pages 311–318.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of NAACL-HLT*, pages 2523–2544.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. [In-context retrieval-augmented language models](#). *Proceedings of TACL*, pages 1316–1331.
- Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, QiaoQiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. [RocketQAv2: A joint training method for dense passage retrieval and passage re-ranking](#). In *Proceedings of EMNLP*, pages 2825–2835.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. [The probabilistic relevance framework: Bm25 and beyond](#). *Foundations and Trends® in Information Retrieval*, (4):333–389.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. [Laion-400m: Open dataset of clip-filtered 400 million image-text pairs](#). *ArXiv preprint*.
- Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. 2019. [KVQA: knowledge-aware visual question answering](#). In *Proceedings of AAAI*, pages 8876–8884.
- Sahel Sharifymoghaddam, Shivani Upadhyay, Wenhu Chen, and Jimmy Lin. 2024. [Unirag: Universal retrieval augmentation for multi-modal large language models](#). *ArXiv preprint*.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. [REPLUG: Retrieval-augmented black-box language models](#). In *Proceedings of NAACL-HLT*, pages 8371–8384.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. [Retrieval augmentation reduces hallucination in conversation](#). In *Proceedings of EMNLP Findings*, pages 3784–3803.
- Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024a. [Generative multimodal models are in-context learners](#). In *Proceedings of CVPR*, pages 14398–14409.
- Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024b. [Emu: Generative pretraining in multimodality](#). In *Proceedings of ICLR*.
- Sahar Tahmasebi, Eric Müller-Budack, and Ralph Ew-erth. 2024. [Multimodal misinformation detection using large vision-language models](#). In *Proceedings of CIKM*, pages 2189–2199.

- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. [Gemini: a family of highly capable multimodal models](#). *ArXiv preprint*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a large-scale dataset for fact extraction and VERification](#). In *Proceedings of NAACL-HLT*, pages 809–819.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. [Llama: Open and efficient foundation language models](#). *ArXiv preprint*.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. [Cider: Consensus-based image description evaluation](#). In *Proceedings of CVPR*, pages 4566–4575.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *ArXiv preprint*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. [Emergent abilities of large language models](#). *Transactions on Machine Learning Research*.
- Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muenighoff, Defu Lian, and Jian-Yun Nie. 2024. [C-pack: Packed resources for general chinese embeddings](#). In *Proceedings of SIGIR*, pages 641–649.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021a. [Approximate nearest neighbor negative contrastive learning for dense text retrieval](#). In *Proceedings of ICLR*.
- Wenhan Xiong, Xiang Lorraine Li, Srini Iyer, Jingfei Du, Patrick S. H. Lewis, William Yang Wang, Yashar Mehdad, Scott Yih, Sebastian Riedel, Douwe Kiela, and Barlas Oguz. 2021b. [Answering complex open-domain questions with multi-hop dense retrieval](#). In *Proceedings of ICLR*.
- Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. [Corrective retrieval augmented generation](#). *ArXiv preprint*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). In *Proceedings of ICLR*.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. [Minicpm-v: A gpt-4v level mllm on your phone](#). *ArXiv preprint*.
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Richard James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-Tau Yih. 2023. [Retrieval-augmented multimodal language modeling](#). In *Proceedings of ICML*, pages 39755–39769.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. [From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions](#). *Proceedings of TACL*, pages 67–78.
- Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. 2023a. [Scaling autoregressive multi-modal models: Pretraining and instruction tuning](#). *ArXiv preprint*.
- Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu. 2023b. [Augmentation-adapted retriever improves generalization of language models as generic plug-in](#). In *Proceedings of ACL*, pages 2421–2436.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. [A survey of large language models](#). *ArXiv preprint*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024. [LlamaFactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of ACL*, pages 400–410.
- Junjie Zhou, Zheng Liu, Shitao Xiao, Bo Zhao, and Yongping Xiong. 2024a. [VISTA: Visualized text embedding for universal multi-modal retrieval](#). In *Proceedings of ACL*, pages 3185–3200.
- Tianshuo Zhou, Sen Mei, Xinze Li, Zhenghao Liu, Chenyan Xiong, Zhiyuan Liu, Yu Gu, and Ge Yu. 2024b. [MARVEL: Unlocking the multi-modal capability of dense retrieval via visual module plugin](#). In *Proceedings of ACL*, pages 14608–14624.

Task	Source	#Docs		#Query	
		Image	Text	Train	Test
Img Cap.	WebQA	383,750	-	3,000	3,000
MM QA.	WebQA	389,750	787,697	3,000	3,000
MM FV.	Factify	41,000	41,000	3,000	3,000
Img Rerank.	WebQA	389,750	-	-	3,000

Table 4: Data Statistics of M²RAG.

A Appendix

A.1 License

We show the licenses of the datasets that we use. WebQA uses CC0-1.0 license, while Factify uses MIT license. All these licenses and agreements permit the use of their data for academic purposes.

A.2 More Details of M²RAG Benchmark

In this section, we provide additional details regarding the data used in the M²RAG benchmark. The data statistics are shown in Table 4.

The M²RAG benchmark incorporates tasks, such as image captioning, multi-modal question answering, and image reranking, all of which are built upon the WebQA (Chang et al., 2022) benchmark. For both the multi-modal question answering and image reranking tasks, the same multi-modal retrieval corpus is used. For the image captioning task, we filter out any image documents that are selected for the construction of both the training and evaluation sets. The image-caption pairs used in the image reranking test set are the same as those in the image captioning test set.

For the multi-modal fact verification task, since the test set labels are not publicly available in the Factify (Mishra et al., 2022) dataset, we follow the approach of Tahmasebi et al. (2024) and sample data from the validation set to create the evaluation set. All text and image documents from the training and validation sets of the Factify dataset are then collected to construct the retrieval corpus. The original Factify dataset consists of five categories: “Support_Text”, “Support_Multimodal”, “Insufficient_Text”, “Insufficient_Multimodal”, and “Refute”. When constructing the training and evaluation datasets for M²RAG, we select an equal number of samples from each of these five categories. Since our RAG scenario involves both text and image information, we consolidate these modality-based categories into three: “Support”, “Refute”, and “Insufficient”, in order to better evaluate the effectiveness of LLMs in the multi-modal fact verification task.

Benchmark	BLEU-2	BLEU-4	ROUGE-L	CIDEr
MiniCPM-V 2.6 (8B)				
MSCOCO	11.53	4.19	30.68	33.48
M ² RAG	4.62	1.91	17.58	18.39
Qwen2-VL (7B)				
MSCOCO	16.93	6.75	37.51	70.75
M ² RAG	5.08	2.24	19.48	26.01

Table 5: Performance of Different MLLMs in the Image Captioning Tasks of MSCOCO and M²RAG.

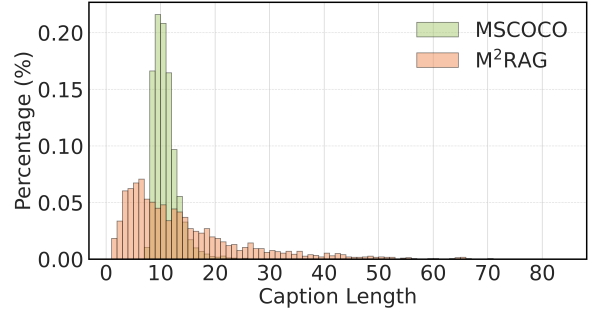


Figure 5: Length Distribution of Captions in the MSCOCO and M²RAG Benchmarks.

A.3 Prompt Templates Used in M²RAG

In Figure 6, we present the prompt templates designed for various task scenarios in M²RAG. Additionally, we describe the input format for the image reranking task. In terms of image placement, we use the placeholder *{image}* for the image, following the method proposed by Hu et al. (2024).

A.4 Comparison of Different Image Captioning Tasks

In this section, we compare the performance of MiniCPM-V 2.6 and Qwen2-VL on the image captioning task using the MSCOCO (Lin et al., 2014) and M²RAG datasets. For the MSCOCO dataset, we use the version employed in the image captioning task of UniRAG (Sharifymoghaddam et al., 2024) and follow the same processing method described in their paper.

As shown in Table 5, both MiniCPM-V 2.6 and Qwen2-VL exhibit lower performance on M²RAG compared to MSCOCO, with average declines of over 9% and 19%, respectively. This suggests that the image captioning task in M²RAG is more challenging for multi-modal large language models (MLLMs) than that in MSCOCO. As shown in Figure 5, the length of image captions in M²RAG is not fixed and varies depending on the scene, with diverse entities and detailed scene descriptions. This indicates the complexity of M²RAG that requires MLLMs to leverage external knowledge to generate

accurate captions. In contrast, the image captions in MSCOCO are more straightforward and regular, enabling MLLMs to perform better by relying primarily on internal knowledge. These observations further underscore the need for M²RAG to serve as a more challenging benchmark for multi-modal RAG tasks.

A.5 Additional Case Studies

As shown in Figure 7, we present additional case studies from various tasks to assess the effectiveness of vanilla RAG MM-RAIT models in utilizing multi-modal context. In the RAG setting, we use the top-5 retrieved multi-modal documents for inference.

In the image captioning task, MLLMs initially provide generic descriptions based on internal knowledge. However, when multi-modal context is incorporated, they are able to extract richer and more specific information, such as identifying landmarks like the “Library of Congress”. After MM-RAIT training, both MiniCPM-V 2.6 and Qwen2-VL generate more accurate captions. A similar improvement is observed in the image reranking task, where vanilla MLLMs initially struggle to align the semantics of the image and caption. After MM-RAIT training, fine-grained alignments between images and captions are achieved, allowing Qwen2-VL to rank the image of “Bellagio and Caesars Palace from the left side of Bellagio’s fountains” first, even though the reranking task is not involved during training.

In multi-modal question answering, the retrieved image, along with the caption “1929 Cadillac Sports Phaeton”, aids Qwen2-VL in producing accurate answers within the RAG framework. However, MiniCPM-V initially struggles to provide a correct answer. After MM-RAIT training, both models consistently provide stable and accurate answers. Similarly, in multi-modal fact verification, vanilla RAG models struggle to extract useful information from noisy documents while MM-RAIT enables the models to better extract and utilize relevant evidence, thereby improving their fact verification performance.

Prompts of Different Tasks w/o Retrieval	Prompts of Different Tasks w/ Retrieval
<p>[Task]: Image Captioning [Instruction]: You are an intelligent assistant capable of generating accurate and detailed captions for images. I will provide you with an image. Response the caption for the image directly, do not include any explanations or unrelated details. <i>{image}</i></p> <p>Caption:</p>	<p>[Task]: Image Captioning [Instruction]: You are an intelligent assistant capable of generating accurate and detailed captions for images. You will be given one image and several retrieved images and their captions. The first image is the input image, others are retrieved examples to help you. Response the caption for the image directly, do not include any explanations or unrelated details. <i>{image}</i> ... <i>{image}</i> Image: The first image. Retrieved_Image_Caption: <i>{retrieval_image_caption}</i></p> <p>Caption:</p>
<p>[Task]: Multi-Modal QA [Instruction]: You are an intelligent assistant capable of answering complex questions. Your task is to carefully analyze the question and provide a detailed, well-structured, and contextually accurate answer in the form of a complete sentence. Write a detailed and accurate answer directly, do not include unrelated details in your response. Question: <i>{question}</i></p> <p>Answer:</p>	<p>[Task]: Multi-Modal QA [Instruction]: You are an intelligent assistant capable of answering complex questions using both visual and textual data. I will provide you with a question and several retrieved images or texts to assist you. Your task is to carefully analyze the question and retrieved information to provide a detailed, well-structured, and contextually accurate answer in the form of a complete sentence. Write a detailed and accurate answer directly, do not include additional context, or unrelated details in your response. <i>{image}</i> ... <i>{image}</i> Question: <i>{question}</i> Retrieved_Image_Caption: <i>{retrieval_image_caption}</i> Retrieved_Text: <i>{retrieval_text}</i></p> <p>Answer:</p>
<p>[Task]: Multi-Modal Fact Verification [Instruction]: You are an intelligent assistant capable of verifying the factual accuracy by classifying data samples into one of three categories: Support, Insufficient, and Refute. I will provide you with a claim (to be verified) consisting of text and an image. Your task is to verify the claim based on your own knowledge into one of the three categories. Response to the category of the claim directly, do not say any other words or explain. <i>{image}</i> Claim_Image: The first image. Claim_Text: <i>{claim_text}</i></p> <p>Category:</p>	<p>[Task]: Multi-Modal Fact Verification [Instruction]: You are an intelligent assistant capable of verifying the factual accuracy by classifying data samples into one of three categories: Support, Insufficient, and Refute. I will provide you with a claim (to be verified) consisting of text and an image, and retrieved images or texts. The first image is the claim_image, and others are retrieved examples to help you. Your task is to classify the relationship between the claim and the evidence from the retrieved documents and your own knowledge into one of the three categories. Response the category for the claim and document directly, do not say any other word or explain. <i>{image}</i> ... <i>{image}</i> Claim_Image: The first image. Claim_Text: <i>{claim_text}</i> Retrieved_Document_Text: <i>{retrieval_text}</i></p> <p>Category:</p>
Input format of Image Reranking Task	
<p>[Task]: Image Reranking [Input Format]: <Rule>: User <Content>: <i>{retrieved_image}</i></p> <p><Role>: Assistant <Content>: <i>{caption}</i></p>	

Figure 6: Prompts Used for Different Tasks in Our M²RAG Benchmark.

Image: 

Ground Truth Caption: Library of Congress ceiling columns Washington DC

Retrieved Documents:





MiniCPM-V 2.6

Vanilla Model: Looking up at the grandiose interior of a historical building, showcasing ornate **columns** and intricately designed **ceilings** with stained glass windows.

Vanilla RAG: Magnificent interior of the **Library of Congress**, showcasing intricate architectural details and vibrant **ceiling** artwork.

MM-RAIT w/o RAG: Looking Up at the Grand Stained Glass **Ceiling** of a **Library**.

MM-RAIT: **Library of Congress - Washington, D.C., 2014 Library of Congress interior in 2014**

Qwen2-VL

Vanilla Model: The interior of a grand building with ornate **columns** and a detailed **ceiling**.

Vanilla RAG: **The Library of Congress in Washington, D.C.**

MM-RAIT w/o RAG: **Library of Congress - interior**


MM-RAIT: **Library of Congress interior, Washington, D.C., 2012 Library of Congress in Washington, D.C. in 2012**

(a) Image Captioning Task.

Question: Does the 1929 Cadillac Sports Phaeton have an extra wheel on its design?

Ground Truth Answer: **Yes**, the 1929 Cadillac Sports Phaeton **has an extra wheel on its side** that does not touch the road.

Retrieved Documents:



Wheels were available in hickory artillery style, wire wheels or solid pressed ...

The engine displacement of the 1930 LaSalle and the 1928-1929 Cadillac Series 341 were ...

MiniCPM-V 2.6

Vanilla Model: The 1929 Cadillac Phaeton **does not have an extra wheel** in its design; it is a standard four-wheel vehicle.

Vanilla RAG: The 1929 Cadillac Sports Phaeton **does not have an extra wheel** on its design. The images and ...

MM-RAIT w/o RAG: **No**, the 1929 Cadillac Sports Phaeton **does not have an extra wheel** on its design.

MM-RAIT Model: The 1929 Cadillac Sports Phaeton **has an extra wheel** on its design.

Qwen2-VL

Vanilla Model: The 1929 Cadillac Sports Phaeton **does not have an extra wheel** on its design; it is a classic car known for its ...

Vanilla RAG: **Yes**, the 1929 Cadillac Sports Phaeton **has an extra wheel** on its design.


MM-RAIT w/o RAG: The 1929 Cadillac Sports Phaeton **has an extra wheel** on its design.

MM-RAIT Model: **Yes**, the 1929 Cadillac Sports Phaeton **has an extra wheel** on its ...

(b) Multi-Modal Question Answering Task.



Claim: Kerala has decided to ban Jio's Internet services from 2021. The government will start its own network services and mobile phones at half of Jio's prices.

Category: **Refute**



Retrieved Documents:

... The Kerala government has neither banned Jio nor started a state-owned Internet ...

MiniCPM-V 2.6

Vanilla Model: **Insufficient**

Vanilla RAG: **Insufficient**

MM-RAIT w/o RAG: **Insufficient**

MM-RAIT: **Refute**

Qwen2-VL

Vanilla Model: **Support**


Vanilla RAG: **Insufficient**

MM-RAIT w/o RAG: **Insufficient**



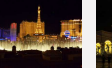


MM-RAIT: **Refute**

(c) Multi-Modal Fact Verification Task.

Caption: Bellagio Caesars Palace Night August 2005 A view of Bellagio and Caesars Palace from the left side of Bellagio's fountains, at night (10 PM local time) in Las Vegas, Nevada, USA.

Ground Truth Image: 

Retrieved Order:

MiniCPM-V 2.6

Vanilla Model:

Image 3 < **Image 4** < **Image 2** < **Image 5** < **Image 1**

11.4375 12.0 12.75 12.75 13.375

MM-RAIT:

Image 5 < **Image 1** < **Image 3** < **Image 2** < **Image 4**

9.625 9.9375 9.9375 10.125 10.25

Qwen2-VL

Vanilla Model:

Image 5 < **Image 1** < **Image 2** < **Image 3** < **Image 4**

12.9375 13.375 14.9375 15.375 16.375

MM-RAIT:

Image 1 < **Image 5** < **Image 2** < **Image 3** < **Image 4**

7.875 8.125 9.1875 9.625 10.4375

(d) Image Reranking Task.

Figure 7: Cases in Different Tasks. For generation tasks, we present the responses of different models using different RAG strategies (w/ or w/o RAG). We use green boxes to mark **the documents that can provide information for the question**. In the model output part, **correct answers** are marked in green, and red for **incorrect**. For Image Reranking task, we presented the order reranked by different models through corresponding PPL scores.