

# MRAG-BENCH: VISION-CENTRIC EVALUATION FOR RETRIEVAL-AUGMENTED MULTIMODAL MODELS

Wenbo Hu<sup>1</sup>, Jia-Chen Gu<sup>1</sup>, Zi-Yi Dou<sup>1</sup>, Mohsen Fayyaz<sup>1</sup>, Pan Lu<sup>2</sup>,  
Kai-Wei Chang<sup>1</sup>, Nanyun Peng<sup>1</sup>

<sup>1</sup>UCLA, <sup>2</sup>Stanford University

{wenbohu, gujc, zdou}@ucla.edu

<https://mragbench.github.io>

## ABSTRACT

Existing multimodal retrieval benchmarks primarily focus on evaluating whether models can retrieve and utilize external *textual knowledge* for question answering. However, there are scenarios where retrieving visual information is either more beneficial or easier to access than textual data. In this paper, we introduce a **multimodal retrieval-augmented generation benchmark**, MRAG-BENCH, in which we systematically identify and categorize scenarios where visually augmented knowledge is better than textual knowledge, for instance, more images from varying viewpoints. MRAG-BENCH consists of 16,130 images and 1,353 human-annotated multiple-choice questions across 9 distinct scenarios. With MRAG-BENCH, we conduct an evaluation of 10 open-source and 4 proprietary large vision-language models (LVLMs). Our results show that all LVLMs exhibit greater improvements when augmented with images compared to textual knowledge, confirming that MRAG-BENCH is vision-centric. Additionally, we conduct extensive analysis with MRAG-BENCH, which offers valuable insights into retrieval-augmented LVLMs. Notably, the top-performing model, GPT-4o, faces challenges in effectively leveraging retrieved knowledge, achieving only a 5.82% improvement with ground-truth information, in contrast to a 33.16% improvement observed in human participants. These findings highlight the importance of MRAG-BENCH in encouraging the community to enhance LVLMs’ ability to utilize retrieved visual knowledge more effectively.

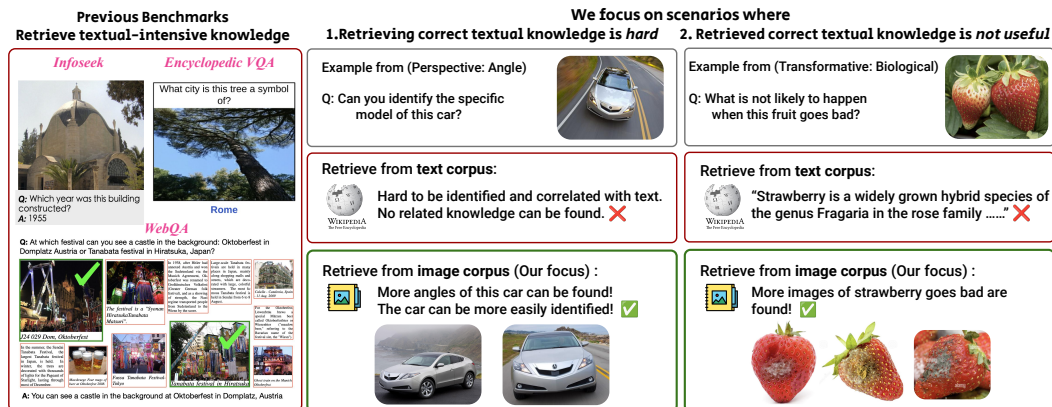


Figure 1: Example scenarios from MRAG-BENCH. Previous benchmarks (Chang et al., 2022; Mensink et al., 2023; Chen et al., 2023b) mainly focused on retrieving from textual knowledge. However, there are scenarios where retrieving correct textual knowledge is hard and sometimes not as useful as visual knowledge.









Benchmarks	Knowledge Modality	Knowledge Source	Multi-Image Input	Diverse Scenarios
K-VQA (Shah et al., 2019)	Text	Wikipedia	✗	✗
OK-VQA (Marino et al., 2019)	Text	Wikipedia	✗	✗
MultiModalQA (Talmor et al., 2021)	Text	Wikipedia	✗	✗
ManyModalQA (Hannan et al., 2020)	Text	Wikipedia	✗	✓
A-OKVQA (Schwenk et al., 2022)	Text	Common/World	✗	✗
ViQuAE (Lerner et al., 2022)	Text	Wikipedia	✗	✗
WebQA (Chang et al., 2022)	Text/Caption	Wikipedia	✗	✗
Encyclopedia VQA (Mensink et al., 2023)	Text	Wikipedia	✗	✗
InfoSeek (Chen et al., 2023b)	Text	Wikipedia	✗	✗
MRAG-BENCH (Ours)	Image	   	✓	✓

Table 1: Compared with previous works, MRAG-BENCH focuses on evaluating LVLMs in utilizing vision-centric retrieval-augmented multimodal knowledge. “Diverse scenarios” refers to whether a benchmark categorized different scenarios during evaluation. : Web, : ImageNet (Russakovsky et al., 2015), : Flowers102 (Nilsback & Zisserman, 2008), : StanfordCars (Krause et al., 2013).

## 1 INTRODUCTION

Retrieval-augmented generation (RAG) has emerged as a promising direction in large vision-language models (LVLMs) (OpenAI, 2023; Liu et al., 2023a; Bai et al., 2023; Huang et al., 2023; Chen et al., 2023a; Hu et al., 2024b; Chen et al., 2024; Tong et al., 2024; McKinzie et al., 2024). By incorporating external knowledge during generation, models such as Wiki-LLaVA (Caffagni et al., 2024) have demonstrated improved performance in knowledge-intensive question answering tasks.

There are several existing benchmarks evaluating retrieval-augmented LVLMs. For example, OK-VQA (Marino et al., 2019) focused on scenarios where the image content alone is insufficient to answer the questions. A-OKVQA (Schwenk et al., 2022) further extended this dataset to incorporate additional types of world knowledge. More recent works (Chang et al., 2022; Chen et al., 2023b; Mensink et al., 2023) further expanded and curated large-scale knowledge base data to evaluate pre-trained vision and language models in knowledge-intensive and information-seeking visual questions. However, as shown in Table 1, these benchmarks remain text-centric, as their questions can often be resolved with related external textual knowledge. In contrast, retrieving visual information is sometimes more beneficial than retrieving text, as humans often gain greater insights from it. Specifically, we illustrate examples in Figure 1 where retrieving correct textual knowledge can be *hard* and retrieved textual knowledge can be *useless*, while retrieving additional images is helpful. For instance, when presented with a top-down view of a car, humans may struggle to accurately identify it; however, with a front-facing view, they can quickly recognize the vehicle and effectively leverage the visual information.

In this paper, we introduce MRAG-BENCH, a benchmark specifically designed for vision-centric evaluation for retrieval-augmented multimodal models, with visual questions typically benefit more from retrieving visual knowledge than textual information. MRAG-BENCH consists of 16,130 images and 1,353 human-annotated multi-choice questions spanning 9 distinctive scenarios. Focusing on utilizing visually augmented knowledge in real-world scenarios, we divide our benchmark into two aspects: *perspective*, where changes in visual entity’s perspective requiring visually augmented knowledge; and *transformative*, where the visual entity undergoes transformative change physically thus requiring visually augmented knowledge. Specifically, MRAG-BENCH requires models to reason about visual entities that undergo perspective changes, such as *angle*, *partial*, *scope* and *occlusion*, as well as transformative changes, such as *temporal*, *incomplete*, *biological* and *deformations*. Additionally, MRAG-BENCH includes 9,673 human-selected images, which serves as the ground-truth image knowledge corpus for model evaluation.

We conduct extensive experiments on MRAG-BENCH to evaluate 10 open-source and 4 proprietary LVLMs. The results confirm that MRAG-BENCH is vision-centric, as all LVLMs show greater improvements when augmented with images compared to textual knowledge. Our results indicate that the best-performing GPT-4o model only achieve 68.68% and 74.5% of accuracy without RAG knowledge and with ground-truth (GT) RAG knowledge, respectively. This substantially outperforms the best open-source model LLaVA-OneVision by 15.39% and 15.52%, respectively.

Statistic	Number
Total questions	1,353
- Multiple-choice questions	1,353 (100%)
- Questions newly annotated	1,353 (100%)
Total Scenarios	9
Unique number of questions	375
Unique number of answers	663
Total number of images	16,130
Unique number of images	16,130
Human selected images	9,673
Average image size (px)	1076 x 851
Maximum question length	20
Maximum answer length	9
Average question length	8.03
Average answer length	2.16
Average choice number	4

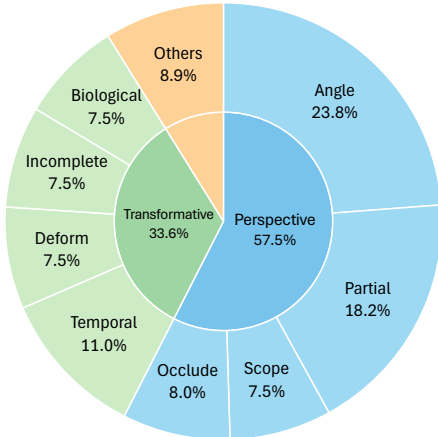


Table 2: Key statistics of MRAG-BENCH. Figure 2: Scenarios distribution of MRAG-BENCH.

Notably, we observe while all models improve with GT knowledge, only proprietary models are able to effectively utilize noisy retrieved multimodal knowledge. This indicates the gap between open-source and close-source models still exists. Open-source models are falling short on their parametric knowledge and the ability to distinguish between high-quality and poor-quality retrieved visually augmented examples. In comparison to humans, GPT-4o achieves only a 5.82% improvement when augmented with GT knowledge and 0.28% with retrieved knowledge, whereas humans demonstrate a 33.16% and 22.91% improvement, respectively. These results highlight the importance of MRAG-BENCH in encouraging the community to develop LVLMs better utilizing of visually augmented knowledge.

## 2 MRAG-BENCH

### 2.1 BENCHMARK OVERVIEW

Our benchmark is designed for systematic evaluation of LVLMs’s vision-centric multimodal RAG abilities. To achieve this, we focus on evaluating the model’s understanding of image objects that are not commonly associated with its knowledge base, while the collected ground-truth images can help incentivize specific visual concepts within LVLMs’ memory. Therefore, we divide our benchmark into two main aspects, as illustrated in the examples in Figure 1:

- *perspective*, refers to the challenges in visual recognition and reasoning that arise when a visual entity is presented from varying viewpoints, scopes, or levels of visibility.
- *transformative*, refers to the challenges that arise when a visual entity undergoes fine-grained physical transformations, making it unfamiliar or not easily associated with the model’s prior knowledge.

MRAG-BENCH consists of 16,130 images and 1,353 multiple choice questions, with key statistics shown in Table 2. MRAG-BENCH adheres to the following design principles: (1) it focuses on real-world scenarios where visually augmented information is useful; (2) it incorporates 9 diverse multimodal RAG scenarios covering various types of image objects; (3) it features cleaned ground-truth images for each question that align with human knowledge; and (4) it provides robust evaluation settings for deterministic evaluations. Unlike previous works focus on retrieving textual knowledge, evaluation on MRAG-BENCH focuses on retrieving vision-centric knowledge, which can be formulated as follows: Given a query tuple  $\mathbf{Q}$  composed of (query image, textual question), the multimodal retriever  $\mathcal{R}$  returns a set of relevant images  $\mathbf{I}$  ( $[\mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_N]$ ), then the LVLm  $\mathcal{M}$  take the input  $(\mathbf{Q}, \mathbf{I})$  and output the final answer.



Figure 3: Qualitative examples on MRAG-BENCH. For each scenario, we show the result of GPT-4o (OpenAI, 2023), Gemini Pro (Team et al., 2023), LLaVA-Next-Interleave (Li et al., 2024b) and Mantis-8B-Siglip (Jiang et al., 2024a). The ground-truth answer is in blue.

## 2.2 BENCHMARK COMPOSITION

MRAG-BENCH provides a systematic evaluation across 9 distinctive multimodal RAG scenarios, with four scenarios focused on the *perspective* understanding of visual entities, four on *transformative* understanding, and one categorized as “others”. As illustrated in Figure 2, each scenario comprises 7.5% to 23.8% of the whole benchmark. The selected examples of each scenario is shown in Figure 3. The details of each scenario are introduced as follows.

**Perspective understanding aspect.** First, we have *perspective* aspect comprising [ANGLE], [PARTIAL], [SCOPE], and [OCCLUSION] dimensions.

- [ANGLE] evaluates the ability of models to utilize visual knowledge of common shooting angles to identify and reason about less common, long-tailed viewpoints of visual entities.
- [PARTIAL] evaluates the ability of models to use complete appearance knowledge to identify and reason when only a partial image of the visual entities is available.
- [SCOPE] evaluates the ability of models to leverage high-resolution, detailed images for identifying and reasoning about visual entities in longer-scoped, low-resolution images.
- [OCCLUSION] evaluates the ability of models to use ground-truth image knowledge to identify and reason when visual entities are occluded or partially hidden in natural scenes.

**Transformative understanding aspect.** On the other hand, the *transformative* understanding scenarios cover [TEMPORAL], [DEFORMATION], [INCOMPLETE], and [BIOLOGICAL] dimensions.

- [TEMPORAL] evaluates the ability of models to use familiar image knowledge to identify and reason about visual entities undergoing temporal changes that may not be represented in the model’s knowledge base.
- [DEFORMATION] evaluates the ability of models to use intact physical appearance knowledge to identify and reason when visual entities undergo deformation not captured in the model’s knowledge base.
- [INCOMPLETE] evaluates the ability of models to compare and contrast the complete layout and structure of image knowledge to identify and reason about missing parts and the correct layout of visual entities.
- [BIOLOGICAL] evaluates the ability of models to utilize image knowledge after biological transformations of the visual entities.

[OTHERS] aims to evaluate the ability of models to leverage geographic image knowledge to accurately identify and reason about the correct regions of origin for the visual entities of interest. All these scenarios work in tandem to comprehensively evaluate LVLMs’ abilities of leveraging visually augmented knowledge.

### 2.3 DATA COLLECTION

As the guidelines discussed in § 2.1, our benchmark collection involves a clean ground-truth image corpus that can resonate with model’s internal knowledge and a query question and image that challenge model’s memory according to our definition of 9 diverse scenarios. To collect a dataset for systematic evaluation of vision-centric multimodal RAG scenarios, we manually annotate all multiple-choice question answering (MCQA) data while sourcing images from either publicly available datasets or manually scraping them from the web.

**Collection of *perspective* aspect.** To collect diverse image objects and knowledge that are not extensively represented in LVLMs’ memories (Zhang et al., 2024c), we considered three sources of data, ImageNet (Russakovsky et al., 2015), Oxford Flowers102 (Nilsback & Zisserman, 2008), and StanfordCars (Krause et al., 2013). To construct a high quality image corpus, for each of the image class that we included in our benchmark, we examined the validation set and excluded the unqualified images which can’t provide sufficient visual information for the recognition of this class. Among the selected corpus, we further humanly picked five representative examples covering the diverse aspects of each class object, as the five ground-truth examples in our experimental results (See §3). For constructing the query images, we adhered to our scenario definitions and manually selected qualified images for the [ANGLE], [SCOPE], and [OCCLUSION] scenarios. For the [PARTIAL] scenario, we randomly cropped images by 50% in both height and width. Then we performed another human inspection to ensure the quality of the cropped images, filtering out examples where the visual object did not occupy the dominant area of the image. We repeated the random cropping process until satisfactory images were obtained, filtering to 20.4 GT images per question on average.

**Collection of *transformative* aspect.** We chose to manually scrape images from the web based on the definitions of the *transformative* aspect. To construct the image corpus, we employed Bing Image Search for each of the image object keyword predefined by us, please refer to Appendix A.1 for more details. We filtered out image objects that did not form a clear transformative pair between the query image and the ground-truth image, retaining approximately 74% of the keyword names in the process. For ground-truth image examples, we employed automatic scripts to download the top 15 images related to its keyword names and human filtered out the unqualified image. On average, this results to 5.9 images per question and the five ground-truth images used during our evaluation are manually selected same as in *perspective* aspect.

According to our guidelines, additional related image object knowledge from the same geographic region can assist in identifying that region more effectively. For the [OTHERS] scenario, we source the data from the GeoDE dataset (Ramaswamy et al., 2023). For each distinct image object category,

Model	Overall	Perspective				Transformative				Others
		Angle	Partial	Scope	Occlusion	Temporal	Deformation	Incomplete	Biological	
Random chance	24.83	27.64	23.98	24.51	19.44	22.15	25.49	29.41	25.49	22.5
Human performance	38.47	25.16	34.96	31.37	41.67	21.48	24.51	58.82	54.9	53.33
+ Retrieved RAG	61.38 <sup>+22.91</sup>	62.42 <sup>+37.26</sup>	60.16 <sup>+35.2</sup>	58.82 <sup>+27.45</sup>	62.96 <sup>+21.29</sup>	54.36 <sup>+32.88</sup>	49.02 <sup>+24.51</sup>	78.43 <sup>+19.61</sup>	63.73 <sup>+8.83</sup>	62.5 <sup>+9.17</sup>
+ GT RAG	71.63 <sup>+33.16</sup>	83.85 <sup>+58.69</sup>	70.33 <sup>+35.37</sup>	66.67 <sup>+35.3</sup>	69.44 <sup>+27.77</sup>	59.73 <sup>+38.25</sup>	68.63 <sup>+44.12</sup>	83.33 <sup>+24.51</sup>	73.53 <sup>+18.63</sup>	69.17 <sup>+15.84</sup>
<i>Open-Source LLMs</i>										
OpenFlamingo-v2-9B	26.83	27.95	26.02	31.37	30.56	29.53	34.31	20.59	17.65	21.67
+ Retrieved RAG	28.31 <sup>+1.48</sup>	29.5 <sup>+1.55</sup>	28.86 <sup>+2.84</sup>	28.43 <sup>-2.94</sup>	30.56 <sup>+0.0</sup>	34.23 <sup>+4.7</sup>	31.37 <sup>-2.94</sup>	22.55 <sup>+1.96</sup>	21.57 <sup>+3.92</sup>	22.5 <sup>+0.83</sup>
+ GT RAG	28.90 <sup>+2.07</sup>	26.71 <sup>-1.24</sup>	33.74 <sup>+7.72</sup>	28.43 <sup>-2.94</sup>	33.33 <sup>+2.77</sup>	35.57 <sup>+6.04</sup>	27.45 <sup>-6.86</sup>	27.45 <sup>+6.86</sup>	25.49 <sup>+7.84</sup>	18.33 <sup>-3.34</sup>
Idefics2-8B	31.04	31.06	33.33	31.37	38.89	30.2	35.29	25.49	24.51	26.67
+ Retrieved RAG	30.16 <sup>-0.88</sup>	29.81 <sup>-1.25</sup>	27.64 <sup>-5.69</sup>	29.41 <sup>-1.96</sup>	36.11 <sup>-2.78</sup>	36.24 <sup>+6.04</sup>	28.43 <sup>+6.86</sup>	27.45 <sup>+1.96</sup>	32.35 <sup>+7.84</sup>	25.83 <sup>-0.84</sup>
+ GT RAG	37.03 <sup>+5.99</sup>	36.34 <sup>+5.28</sup>	35.37 <sup>+2.04</sup>	38.24 <sup>+6.87</sup>	54.63 <sup>+15.74</sup>	47.65 <sup>+17.45</sup>	36.27 <sup>+0.98</sup>	24.51 <sup>-0.98</sup>	34.31 <sup>+9.8</sup>	25.83 <sup>-0.84</sup>
VILA1.5-13B	43.68	45.34	41.87	52.94	48.15	50.34	38.24	21.57	30.39	57.5
+ Retrieved RAG	35.48 <sup>-8.2</sup>	33.54 <sup>-11.8</sup>	28.86 <sup>-13.01</sup>	29.41 <sup>-23.53</sup>	40.74 <sup>-7.41</sup>	47.65 <sup>-2.69</sup>	33.33 <sup>+4.91</sup>	22.55 <sup>-0.98</sup>	33.33 <sup>+2.94</sup>	54.17 <sup>-3.33</sup>
+ GT RAG	47.01 <sup>+3.33</sup>	45.65 <sup>+0.31</sup>	46.75 <sup>+8.88</sup>	39.22 <sup>-13.72</sup>	51.85 <sup>+3.7</sup>	53.69 <sup>+3.35</sup>	43.14 <sup>+4.9</sup>	25.49 <sup>+3.92</sup>	44.12 <sup>+13.73</sup>	69.17 <sup>+11.67</sup>
Mantis-8B-clip-llama3	40.8	45.03	39.43	42.16	49.07	49.66	36.27	28.43	19.61	45.0
+ Retrieved RAG	36.88 <sup>-3.92</sup>	36.65 <sup>-8.38</sup>	34.96 <sup>-4.47</sup>	42.16 <sup>0.0</sup>	47.22 <sup>-1.85</sup>	50.34 <sup>+0.68</sup>	33.33 <sup>-2.94</sup>	18.63 <sup>-9.8</sup>	21.57 <sup>+1.96</sup>	42.5 <sup>-2.5</sup>
+ GT RAG	44.72 <sup>+3.92</sup>	48.14 <sup>+3.11</sup>	46.75 <sup>+7.32</sup>	43.14 <sup>+0.98</sup>	54.63 <sup>+5.56</sup>	57.05 <sup>+7.39</sup>	45.1 <sup>+8.83</sup>	19.61 <sup>-8.82</sup>	18.63 <sup>-0.98</sup>	51.67 <sup>+6.67</sup>
Mantis-8B-siglip-llama3	45.01	46.89	45.12	57.84	58.33	45.64	45.1	26.47	29.41	45.0
+ Retrieved RAG	39.62 <sup>-5.39</sup>	42.55 <sup>-4.34</sup>	35.37 <sup>-9.75</sup>	47.06 <sup>-10.78</sup>	47.22 <sup>-11.11</sup>	42.95 <sup>-2.69</sup>	45.1 <sup>0.0</sup>	23.53 <sup>-2.94</sup>	29.41 <sup>0.0</sup>	40.83 <sup>-4.17</sup>
+ GT RAG	48.85 <sup>+3.84</sup>	54.66 <sup>+7.77</sup>	52.85 <sup>+7.73</sup>	51.96 <sup>-5.88</sup>	58.33 <sup>0.0</sup>	48.99 <sup>+3.35</sup>	50.0 <sup>+4.9</sup>	21.57 <sup>-4.9</sup>	33.33 <sup>+3.92</sup>	49.17 <sup>+1.67</sup>
Deepseek-VL-7B-chat	43.39	45.34	47.56	47.06	45.37	46.31	48.04	28.43	20.59	49.17
+ Retrieved RAG	34.66 <sup>-8.73</sup>	33.54 <sup>-11.8</sup>	32.11 <sup>-15.45</sup>	33.33 <sup>-13.73</sup>	37.04 <sup>-8.33</sup>	43.62 <sup>-2.69</sup>	40.2 <sup>-7.84</sup>	20.59 <sup>-7.84</sup>	26.47 <sup>+5.88</sup>	45.0 <sup>-4.17</sup>
+ GT RAG	50.33 <sup>+6.94</sup>	54.04 <sup>+8.7</sup>	56.5 <sup>+8.94</sup>	50.98 <sup>+3.92</sup>	56.48 <sup>+11.11</sup>	57.05 <sup>+10.74</sup>	50.0 <sup>+1.96</sup>	21.57 <sup>-6.86</sup>	23.53 <sup>+2.94</sup>	60.83 <sup>+11.66</sup>
LLaVA-NeXT-Interleave-7B	43.46	44.41	43.5	40.2	64.81	44.97	44.12	32.35	26.47	45.83
+ Retrieved RAG	40.35 <sup>-3.11</sup>	40.06 <sup>-4.35</sup>	33.33 <sup>-10.17</sup>	39.22 <sup>-0.98</sup>	56.48 <sup>-8.33</sup>	43.62 <sup>-1.35</sup>	44.12 <sup>0.0</sup>	27.45 <sup>-4.9</sup>	36.27 <sup>+9.8</sup>	49.17 <sup>+3.34</sup>
+ GT RAG	52.99 <sup>+9.53</sup>	54.97 <sup>+10.56</sup>	54.88 <sup>+11.38</sup>	49.02 <sup>+8.82</sup>	62.04 <sup>-2.77</sup>	52.35 <sup>+7.38</sup>	47.06 <sup>+2.94</sup>	38.24 <sup>+5.89</sup>	48.04 <sup>+21.57</sup>	61.67 <sup>+15.84</sup>
mPLUG-Owl3-7B	49.74	48.45	50.81	54.9	58.33	54.36	51.96	30.39	45.1	51.67
+ Retrieved RAG	41.83 <sup>-7.91</sup>	40.06 <sup>-8.39</sup>	36.59 <sup>-14.22</sup>	40.2 <sup>-14.7</sup>	50.0 <sup>-8.33</sup>	50.34 <sup>-4.02</sup>	46.08 <sup>-5.88</sup>	20.59 <sup>-9.8</sup>	51.96 <sup>+6.86</sup>	46.67 <sup>-5.0</sup>
+ GT RAG	56.32 <sup>+6.58</sup>	58.39 <sup>+9.94</sup>	58.94 <sup>+8.13</sup>	58.82 <sup>+3.92</sup>	62.96 <sup>+4.63</sup>	61.74 <sup>+7.38</sup>	59.8 <sup>+7.84</sup>	26.47 <sup>-3.92</sup>	50.0 <sup>+4.9</sup>	58.33 <sup>+6.66</sup>
LLaVA-OneVision	53.29	58.39	56.1	49.02	60.19	47.65	53.92	37.25	52.94	51.67
+ Retrieved RAG	50.11 <sup>-3.18</sup>	50.93 <sup>-7.46</sup>	48.78 <sup>-7.32</sup>	50.0 <sup>-0.98</sup>	60.19 <sup>0.0</sup>	50.34 <sup>-2.69</sup>	48.04 <sup>-5.88</sup>	33.33 <sup>-3.92</sup>	53.92 <sup>-0.98</sup>	54.17 <sup>+2.5</sup>
+ GT RAG	58.98 <sup>+5.69</sup>	62.42 <sup>+4.03</sup>	63.82 <sup>+7.72</sup>	59.8 <sup>+10.78</sup>	66.67 <sup>+6.48</sup>	59.73 <sup>+12.08</sup>	53.92 <sup>-0.0</sup>	30.39 <sup>-6.86</sup>	57.84 <sup>+4.9</sup>	60.83 <sup>+9.16</sup>
Pixtral-12B	47.97	52.48	45.53	58.82	50.0	51.68	49.02	38.24	42.16	37.5
+ Retrieved RAG	45.97 <sup>-2.0</sup>	51.86 <sup>-0.62</sup>	40.24 <sup>-5.29</sup>	53.92 <sup>-4.9</sup>	50.93 <sup>+0.93</sup>	49.66 <sup>-2.02</sup>	47.06 <sup>-1.96</sup>	19.61 <sup>-18.63</sup>	47.06 <sup>+4.9</sup>	46.67 <sup>+9.17</sup>
+ GT RAG	59.28 <sup>+11.31</sup>	63.04 <sup>+10.56</sup>	63.41 <sup>+17.88</sup>	65.69 <sup>+6.87</sup>	66.67 <sup>+16.67</sup>	61.74 <sup>+10.06</sup>	59.8 <sup>+10.78</sup>	20.59 <sup>-17.65</sup>	50.98 <sup>+8.82</sup>	65.0 <sup>+27.5</sup>
<i>Proprietary LLMs</i>										
GPT-4-Turbo	57.21	64.29	59.35	54.9	56.48	62.42	47.06	41.18	59.8	50.0
+ Retrieved RAG	58.95 <sup>+1.74</sup>	66.53 <sup>+2.24</sup>	59.94 <sup>+0.59</sup>	53.94 <sup>-0.96</sup>	66.74 <sup>+10.26</sup>	59.73 <sup>-2.69</sup>	49.06 <sup>+2.0</sup>	38.27 <sup>-2.91</sup>	62.78 <sup>+2.98</sup>	58.83 <sup>+8.83</sup>
+ GT RAG	62.85 <sup>+5.64</sup>	68.94 <sup>+4.65</sup>	69.51 <sup>+10.16</sup>	60.78 <sup>+5.88</sup>	67.59 <sup>+11.11</sup>	63.33 <sup>+0.91</sup>	51.96 <sup>+4.9</sup>	38.24 <sup>-2.94</sup>	59.8 <sup>0.0</sup>	62.5 <sup>+12.5</sup>
Gemini Pro	61.71	68.01	69.92	73.53	71.3	70.47	42.16	39.22	53.92	40.83
+ Retrieved RAG	65.93 <sup>+4.22</sup>	73.29 <sup>+5.28</sup>	69.92 <sup>+0.0</sup>	69.61 <sup>-3.92</sup>	73.15 <sup>+1.85</sup>	75.84 <sup>+5.37</sup>	49.02 <sup>+6.86</sup>	34.31 <sup>-4.91</sup>	56.86 <sup>+2.94</sup>	65.0 <sup>+24.17</sup>
+ GT RAG	71.40 <sup>+9.69</sup>	77.33 <sup>+9.32</sup>	79.27 <sup>+9.35</sup>	78.43 <sup>+4.9</sup>	75.93 <sup>+4.63</sup>	78.52 <sup>+8.05</sup>	54.9 <sup>+12.74</sup>	36.27 <sup>-2.95</sup>	61.76 <sup>+7.84</sup>	72.5 <sup>+31.67</sup>
Claude 3.5 Sonnet	59.87	70.19	57.72	56.86	57.41	68.46	48.04	49.02	62.75	47.5
+ Retrieved RAG	63.56 <sup>+3.69</sup>	73.91 <sup>+3.72</sup>	70.73 <sup>+13.01</sup>	56.86 <sup>+0.0</sup>	62.96 <sup>+5.55</sup>	70.47 <sup>+2.01</sup>	55.88 <sup>+7.84</sup>	31.37 <sup>-17.65</sup>	62.75 <sup>+0.0</sup>	53.33 <sup>+5.83</sup>
+ GT RAG	71.10 <sup>+11.23</sup>	78.88 <sup>+8.69</sup>	80.49 <sup>+22.77</sup>	76.47 <sup>+19.61</sup>	70.37 <sup>+12.96</sup>	75.17 <sup>+6.71</sup>	67.65 <sup>-19.61</sup>	36.27 <sup>-12.75</sup>	65.69 <sup>+2.94</sup>	59.17 <sup>+11.67</sup>
GPT-4o	68.68	76.09	70.42	69.61	74.07	73.82	61.21	47.62	58.82	65.83
+ Retrieved RAG	68.96 <sup>+0.28</sup>	77.95 <sup>+1.86</sup>	78.86 <sup>+8.44</sup>	69.61 <sup>+0.0</sup>	75.0 <sup>+0.93</sup>	73.83 <sup>+0.01</sup>	54.9 <sup>+7.28</sup>	26.47 <sup>-34.74</sup>	59.8 <sup>+0.98</sup>	68.33 <sup>+2.5</sup>
+ GT RAG	74.50 <sup>+5.82</sup>	84.47 <sup>+8.38</sup>	77.46 <sup>+7.04</sup>	82.35 <sup>+12.74</sup>	79.63 <sup>+5.56</sup>	77.18 <sup>+3.36</sup>	68.62 <sup>+7.41</sup>	30.95 <sup>-16.67</sup>	62.75 <sup>+3.93</sup>	80.0 <sup>+14.17</sup>

Table 3: Accuracy scores on MRAG-BENCH. The highest scores for open-source models in each section and proprietary models are highlighted in blue and red, respectively. CLIP retriever is consistently used across all models. Both Retrieved RAG and GT RAG employ top-5 image examples (except for the incomplete scenario, where a single example is intuitively sufficient). The relative difference in performance compared to the score without RAG is shown in subscript, with blue indicating performance drops and red indicating improvements.

we randomly sampled 3 out of 6 regions to serve as the answers for each question and selected the corresponding image as the query image.

**Quality control.** After constructing the entire benchmark, we implemented two quality control procedures: an automatic check with predefined rules and a manual examination of each instance. The automatic check verifies the correct MCQA format, assesses image validity and filters out redundant images in the corpus, more details are presented in Appendix A.1. The manual examination is conducted by two experts in this field, who checked the correspondence between query images and ground-truth image examples, and filtered or revised ambiguous questions and uncorrelated query image and ground-truth images.

### 3 EXPERIMENTS

In this section, we first introduce the experimental setup and evaluation metric (§ 3.1). Then, we present a comprehensive evaluation of 14 recent LVLMs (§ 3.2). We demonstrate the importance of visual knowledge and discuss the critical findings revealed by the results from MRAG-BENCH.

#### 3.1 EXPERIMENTAL SETUP

We evaluate 14 popular LVLMs on MRAG-BENCH, including 4 proprietary models and 10 open-sourced models that can accept multi-image inputs:

- **Proprietary models:** GPT-4o (0513) (OpenAI, 2023), GPT-4-Turbo (OpenAI, 2023), Gemini Pro (Team et al., 2023), and Claude 3.5 Sonnet (Anthropic, 2024).
- **Open-source models:** OpenFlamingo (v2-9B) (Awadalla et al., 2023), Idefics (v2-8B) (Laurençon et al., 2024), VILA (v1.5-13B) (Lin et al., 2023), LLaVA-NeXT-Interleave-7B (Li et al., 2024b), LLaVA-OneVision (Li et al., 2024a), Mantis (clip-llama3, and siglip-llama3 versions; 8B) (Jiang et al., 2024a), mPLUG-Owl3-7B (Ye et al., 2024), Deepseek-VL-7B-chat (Lu et al., 2024a), and Pixtral-12B (Team, 2024).

**Evaluation setup.** We follow standard MCQA evaluation setup and employ accuracy score as our metric. We adopt default generation hyper-parameters selected by each model. Following Lu et al. (2024b), we employ GPT-3.5-turbo to extract the multiple choice answer in rare cases where our pre-defined automatic extraction rules failed. We refer the readers to Appendix A.1 and B for more details on evaluation prompts for both without multimodal RAG and with multimodal RAG scenarios, answer extraction prompt and human performance evaluation protocol.

#### 3.2 MAIN RESULTS

As shown in Table 3, the average performance of the most advanced LVLMs is not better than 68.68% without multimodal RAG knowledge, and 74.5% with ground-truth knowledge, which demonstrates MRAG-BENCH to be a challenging benchmark. The mean accuracies of open-source LVLMs are between 26.83% and 53.29% without RAG knowledge and between 28.90% and 59.28% with ground-truth knowledge, which fall behind from advanced proprietary LVLMs. Notably, MRAG-BENCH proves to be knowledge-intensive as average humans achieved 38.47% without RAG knowledge, while proprietary LVLMs generally perform well, suggesting that their extensive training data equips them with a broader knowledge base. However, when provided with either retrieved or ground-truth knowledge, humans achieve the most significant improvements of 22.91% and 33.16%, respectively. This underscores the need of LVLMs to better utilize visually augmented information like humans.

**Can LVLMs utilize retrieved and ground-truth image knowledge well?** As illustrated in Table 3, all models demonstrate improvement when ground-truth image RAG knowledge is provided. Among the open-source models, they achieve improvements ranging from 2.07% to 11.31% when using ground-truth RAG knowledge, whereas 5.64% to 9.69% improvements are observed from proprietary LVLMs. Interestingly, when images from the multimodal retriever is provided, almost all open-source LVLMs on average demonstrate a declined performance while proprietary models can still gain improvement. This indicates proprietary models possess emerging abilities to distinguish between good and bad image knowledge sources, which is a critical skill in the multimodal RAG domain. We further conducted a qualitative analysis to investigate the reasons behind this, as detailed in the following paragraphs.

**Fine-grained results.** We also report fine-grained scores across 9 scenarios on MRAG-BENCH in Table 3. Remarkably, GPT-4o surpasses most other baselines in various categories, with exceptions in problems related to partial, incomplete and biological scenarios. Notably, GPT-4o outperforms human performance on all perspective aspect as well as on temporal and deformation scenarios within the transformative aspect. We conjecture that incomplete and biological scenarios are less likely to be included in the training knowledge. Interestingly, all models exhibit a decline in performance on incomplete scenarios, with only a few exceptions, while humans find this



Figure 4: Qualitative Example of Proprietary model (Gemini Pro) identifies and utilizes correct examples, while open-source model (LLaVA-Next-Interleave) is misled by noisy retrieved information, resulting in incorrect answers.

Model	Overall	Perspective				Transformative				Others
		Angle	Partial	Scope	Occlusion	Temporal	Deformation	Incomplete	Biological	
LLaVA-NeXT-Interleave-7B	43.46	44.41	43.5	40.2	64.81	44.97	44.12	32.35	26.47	45.83
+ Retrieved Text RAG	37.99 <sup>-5.47</sup>	37.58 <sup>-6.83</sup>	34.96 <sup>-8.54</sup>	33.33 <sup>-6.87</sup>	50.0 <sup>-14.81</sup>	41.61 <sup>-3.36</sup>	35.29 <sup>-8.83</sup>	30.39 <sup>-1.96</sup>	27.45 <sup>+0.98</sup>	51.67 <sup>+5.84</sup>
+ Retrieved Image RAG	40.35 <sup>-3.11</sup>	40.06 <sup>-4.35</sup>	33.33 <sup>-10.17</sup>	39.22 <sup>-0.98</sup>	56.48 <sup>-8.33</sup>	43.62 <sup>-1.35</sup>	44.12 <sup>-0.0</sup>	27.45 <sup>-4.9</sup>	36.27 <sup>+9.8</sup>	49.17 <sup>+3.34</sup>
+ GT Text RAG	41.09 <sup>-2.37</sup>	41.93 <sup>-2.48</sup>	39.02 <sup>-4.48</sup>	38.24 <sup>-1.96</sup>	56.48 <sup>-8.33</sup>	44.97 <sup>+0.0</sup>	43.14 <sup>-0.98</sup>	30.39 <sup>-1.96</sup>	21.57 <sup>-4.9</sup>	50.83 <sup>+5.0</sup>
+ GT Text RAG	52.99 <sup>+9.53</sup>	54.97 <sup>+10.56</sup>	54.88 <sup>+11.38</sup>	49.02 <sup>+8.82</sup>	62.04 <sup>-2.77</sup>	52.35 <sup>+7.38</sup>	47.06 <sup>-2.94</sup>	38.24 <sup>+5.89</sup>	48.04 <sup>+21.57</sup>	61.67 <sup>+15.84</sup>
+ GT Image & Text RAG	47.82 <sup>+4.36</sup>	47.83 <sup>+3.42</sup>	48.78 <sup>+5.28</sup>	44.12 <sup>+3.92</sup>	58.33 <sup>-6.48</sup>	49.66 <sup>+4.69</sup>	48.04 <sup>+3.92</sup>	30.39 <sup>-1.96</sup>	35.29 <sup>+8.82</sup>	62.5 <sup>+16.67</sup>
GPT-4-Turbo	57.21	64.29	59.35	54.9	56.48	62.42	47.06	41.18	59.8	50.0
+ Retrieved Text RAG	56.61 <sup>-0.6</sup>	61.8 <sup>-2.49</sup>	59.35 <sup>-0.0</sup>	59.8 <sup>+4.9</sup>	58.33 <sup>+1.85</sup>	59.06 <sup>-3.36</sup>	49.02 <sup>-1.96</sup>	33.33 <sup>-7.85</sup>	60.78 <sup>+0.98</sup>	52.5 <sup>+2.5</sup>
+ Retrieved Image RAG	58.95 <sup>+1.74</sup>	66.53 <sup>+2.24</sup>	59.94 <sup>+0.59</sup>	53.94 <sup>-0.96</sup>	66.74 <sup>+10.26</sup>	59.73 <sup>-2.69</sup>	49.06 <sup>-2.0</sup>	38.27 <sup>-2.91</sup>	62.78 <sup>+2.98</sup>	58.83 <sup>+8.83</sup>
+ GT Text RAG	58.98 <sup>+1.77</sup>	68.01 <sup>+3.72</sup>	63.41 <sup>+4.06</sup>	65.69 <sup>+10.79</sup>	63.89 <sup>+7.41</sup>	59.73 <sup>-2.69</sup>	38.24 <sup>-8.82</sup>	37.25 <sup>-3.93</sup>	58.82 <sup>-0.98</sup>	50.83 <sup>+0.83</sup>
+ GT Image RAG	62.85 <sup>+5.64</sup>	68.94 <sup>+4.65</sup>	69.51 <sup>+10.16</sup>	60.78 <sup>+5.88</sup>	67.59 <sup>+11.11</sup>	63.33 <sup>+0.91</sup>	51.96 <sup>+4.9</sup>	38.24 <sup>-2.94</sup>	59.8 <sup>+0.0</sup>	62.5 <sup>+12.5</sup>
+ GT Image & Text RAG	65.11 <sup>+7.9</sup>	72.05 <sup>+7.76</sup>	72.76 <sup>+13.41</sup>	67.65 <sup>+12.75</sup>	70.37 <sup>+13.89</sup>	71.81 <sup>+9.39</sup>	46.08 <sup>-0.98</sup>	39.22 <sup>-1.96</sup>	60.78 <sup>+0.98</sup>	57.5 <sup>+7.5</sup>

Table 4: LVLMS performance on MRAG-BENCH with textual knowledge v.s visual knowledge. Both the open-source and proprietary model benefit more from image knowledge.

task relatively easy, achieving 58.82% and 83.33% scores with ground-truth knowledge. This further highlights the importance of leveraging retrieved visually augmented knowledge to address questions that do not directly incentivize knowledge stored in the models’ memories.

**Why can proprietary models better utilize retrieved images?** We conduct an error analysis on an open-source model (LLaVA-Next-Interleave) and a proprietary model (Gemini Pro). For a fair comparison, we filtered results where LLaVA-Next-Interleave answered correctly without or with GT knowledge but was misled to wrong answer with retrieved examples. One example is illustrated in Figure 4, the retrieved images contain two correct examples and three false examples. While Gemini Pro is able to utilize all retrieved images, LLaVA-Next-Interleave leverages bad examples and makes wrong prediction. This example helps explain why do almost all open-source models have lower performance with retrieved knowledge.

#### 4 ANALYSIS

In this section, we conduct quantitative analysis addressing three important questions: 1) To what extent can LVLMS benefit more from visual knowledge than from textual knowledge on MRAG-BENCH? (§ 4.1) 2) How does the performance of LVLMS vary with examples retrieved from different retrievers? (§ 4.2) 3) How many ground-truth visual knowledge examples are required for LVLMS to continue benefiting? (§ 4.3)

##### 4.1 HOW MUCH CAN VISUAL KNOWLEDGE BENEFIT MORE THAN TEXTUAL KNOWLEDGE?

We used the Wikipedia corpus as of 2023/07/01 as our text knowledge corpus<sup>1</sup>. To ensure a fair comparison, we employed the same multimodal retriever (CLIP) for retrieving either text or image knowledge. The top-5 ranked documents or images are used for augmenting the input. We selected one open-source (LLaVA-Next-Interleave) and one proprietary (GPT-4-Turbo) LVLMS to examine their preference for textual knowledge versus image knowledge on MRAG-BENCH. As shown

<sup>1</sup><https://www.kaggle.com/datasets/jjinho/wikipedia-20230701>



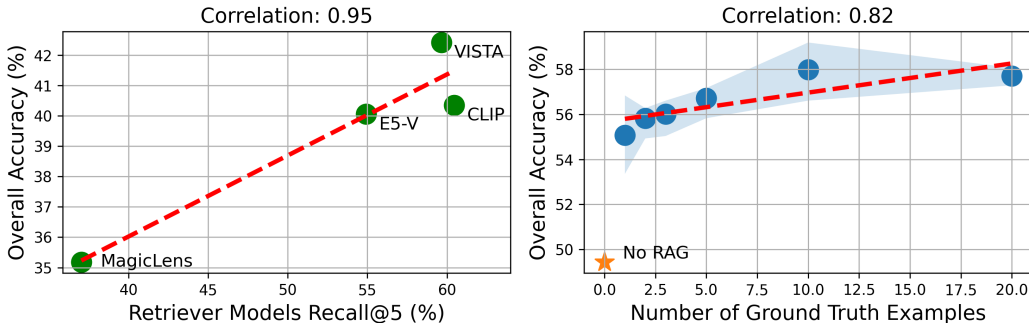


Figure 5: Left: LLaVA-Next-Interleave results with 4 different multimodal retrievers. Its performance using retrieved images correlates 95% with retriever’s Recall@5 scores. Right: Average results of three random seed runs. Improve the number of ground-truth RAG examples shows steady increase of model’s performance, reaches the maximum with 10 examples.

in Table 4, when both models utilized retrieved knowledge, LLaVA-Next-Interleave demonstrated a 2.36% improvement with image knowledge over text knowledge, while GPT-4-Turbo showed a 2.34% improvement. When using GT knowledge, LLaVA-Next-Interleave exhibited an 11.90% improvement with image knowledge over text knowledge, compared to a 3.87% improvement for GPT-4-Turbo. Interestingly, when both GT image and text knowledge are provided, LLaVA-Next-Interleave indicated less improvement than with GT image alone whereas GPT-4-Turbo further pushed its performance. All these results demonstrate that retrieving visual knowledge is more helpful than retrieving text on MRAG-BENCH.

#### 4.2 HOW DOES RETRIEVER PERFORMANCE AFFECT LVLMS?

We picked four recent best-performing multimodal retrievers, including CLIP (Radford et al., 2021), MagicLens (Zhang et al., 2024a), E5-V (Jiang et al., 2024b), VISTA (Zhou et al., 2024) and evaluated their performance (Recall@5). The detailed retriever performance can be found at Table 6 in Appendix C. We selected LLaVA-Next-Interleave as the end model to assess its performance. As shown in Figure 5, when retrievers achieve higher Recall@5 scores (i.e., better retrieved examples), the LLaVA-Next-Interleave’s accuracy tends to improve, demonstrating a strong 95% positive correlation. Interestingly, despite similar Recall@5 scores from CLIP and VISTA retrievers, LLaVA-Next-Interleave demonstrated a 2.07% gap in overall accuracy. We conjecture that the order of the correctly retrieved examples may also impact the model’s final performance. The sensitivity to the order of retrieved examples is a common issue that persists across various models. Although this phenomenon, known as position bias, has been examined in text-based RAG (Lu et al., 2022b; Wang et al., 2023), its impact on visual RAG remains unexplored, presenting a promising direction for future research.

#### 4.3 HOW MANY GROUND-TRUTH IMAGE EXAMPLES ARE NEEDED?

For simplicity, all our experiments used five retrieved or ground-truth image examples. However, it is worth exploring how many examples LVLMS can effectively leverage. As noted in § 2.3, the perspective aspect of our benchmark includes an average of 20.4 ground-truth examples. To investigate further, we perform an analysis focusing on the perspective and others aspects, covering a total of 892 questions. As shown in Figure 5, we evaluated LLaVA-Next-Interleave using 1, 2, 3, 5, 10, 20 GT examples, averaging the results across three random seeds for sampling the GT examples. LLaVA-Next-Interleave saw the greatest improvement of 5.64% with just one GT example. Performance continued to increase steadily, reaching a peak at 10 GT examples, which was 0.29% higher than with 20 GT examples. One possible explanation could be LLaVA-Next-Interleave may not be able to better leverage visually augmented knowledge in long context scenarios. Moreover, the complexity of questions affects the number of images needed too, one ground-truth example sometimes help the model the most on MRAG-BENCH. We encourage the research on adaptatively deciding the number of necessary images based on the complexity of questions.

## 5 RELATED WORK

We overview three lines of related work: 1) multimodal retrieval-augmented generation benchmarks (§ 5.1), 2) large vision language models (§ 5.2), and 3) retrieval-augmented multimodal large language models (§ 5.3).

### 5.1 MULTIMODAL RETRIEVAL-AUGMENTED GENERATION BENCHMARKS

A number of recent benchmarks have been developed to comprehensively assess the capabilities of LVLMs (Lu et al., 2021; 2022a; Li et al., 2023; Liu et al., 2023c; Lu et al., 2024b; Yue et al., 2023; Zhang et al., 2024b; Ying et al., 2024). There are several benchmarks well-suited for evaluating retrieval-augmented LVLMs. For instance, OK-VQA (Marino et al., 2019) and A-OKVQA (Schwenk et al., 2022) both focus on scenarios where external textual knowledge is required to answer visual questions. More recent works (Chang et al., 2022; Chen et al., 2023b; Mensink et al., 2023) have curated large-scale knowledge bases to evaluate models on knowledge-intensive and information-seeking visual questions. In contrast, MRAG-BENCH focus on scenarios where retrieving visual information is more helpful than retrieving text.

### 5.2 LARGE VISION LANGUAGE MODELS

Large Vision Language Models (LVLMs) (Liu et al., 2023b; Zhu et al., 2023a; Dai et al., 2023; Yin et al., 2023; Hu et al., 2024a) have showcased promising results on a wide variety of vision-language tasks. Many works, such as Flamingo (Alayrac et al., 2022), Emu (Sun et al., 2023), Idefics (Laurençon et al., 2023), and VILA (Lin et al., 2023), have demonstrated in-context learning capabilities, where multiple image examples can be leveraged to improve text generation. Recent works start training LVLMs with interleaved image-text corpora, such as MMC4 (Zhu et al., 2023b) and OBELICS (Laurençon et al., 2023), for pretraining, as well as high-quality instruction tuning in models like Mantis-Instruct (Jiang et al., 2024a), LLaVA-Next-Interleave (Li et al., 2024b), and LLaVA-OneVision (Li et al., 2024a), enabling models to process and understand information from multiple images. Naturally, evaluating the ability of LVLMs to effectively leverage visually augmented knowledge becomes an important task, which is the primary focus of MRAG-BENCH.

### 5.3 RETRIEVAL-AUGMENTED MULTIMODAL LARGE LANGUAGE MODELS

Retrieval-augmented generation (RAG) has emerged as a potential solution to overcome limitations in language models by incorporating external knowledge retrieval during the generation process (Lewis et al., 2020; Shi et al., 2024). Reasonably, several works have focused on using multimodal knowledge to enhance the generation capabilities of Large Language Models (LLMs) (Yasunaga et al., 2023; Chen et al., 2022; Zhao et al., 2023; Cui et al., 2024). Recently, more works (Caffagni et al., 2024; Xuan et al., 2024; Du et al., 2024) has incorporated external knowledge to improve LVLMs’ general generation abilities and the comprehensiveness of their reasoning. Although some works (Chen et al., 2022; Yuan et al., 2023) have proposed directly using image information from the web, a systematic vision-centric benchmark to evaluate LVLMs’ abilities to leverage visually augmented knowledge is lacking, which is the focus of our work.

## 6 CONCLUSION

In this work, we introduce MRAG-BENCH, a benchmark specifically designed for vision-centric evaluation for retrieval-augmented multimodal models. Our evaluation of 14 LVLMs highlights that visually augmented knowledge brings more improvements on MRAG-BENCH compared to textual knowledge. Moreover, the top-performing model, GPT-4o, struggles to effectively utilize the retrieved knowledge, achieving only a 5.82% improvement when augmented with relevant information, compared to a 33.16% improvement demonstrated by human participants. We further conduct extensive analysis and propose several promising directions for future research. Our findings underscore the significance of MRAG-BENCH in motivating the community to develop LVLMs that better utilize retrieved visual knowledge.

## ACKNOWLEDGMENTS

We would like to thank the reviewers for their valuable reviews. We would also like to thank UCLA-NLP members for their helpful comments and thank PLUSLab members for their proofreading during the paper clinic session. This research is supported in part by the ECOLE program under Cooperative Agreement HR00112390060 with the US Defense Advanced Research Projects Agency (DARPA), and UCLA-Amazon Science Hub for Humanity and Artificial Intelligence.

## REFERENCES

- 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. Flamingo: a visual language model for few-shot learning. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/960a172bc7fbf0177cccb411a7d800-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/960a172bc7fbf0177cccb411a7d800-Abstract-Conference.html). 10
- Anthropic. Claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024. 7
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models. *ArXiv preprint*, abs/2308.01390, 2023. URL <https://arxiv.org/abs/2308.01390>. 7
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *ArXiv preprint*, abs/2308.12966, 2023. URL <https://arxiv.org/abs/2308.12966>. 2
- Davide Caffagni, Federico Cocchi, Nicholas Moratelli, Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Wiki-llava: Hierarchical retrieval-augmented generation for multimodal llms, 2024. URL <https://arxiv.org/abs/2404.15406>. 2, 10
- Yingshan Chang, Guihong Cao, Mridu Narang, Jianfeng Gao, Hisami Suzuki, and Yonatan Bisk. Webqa: Multihop and multimodal QA. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 16474–16483. IEEE, 2022. doi: 10.1109/CVPR52688.2022.01600. URL <https://doi.org/10.1109/CVPR52688.2022.01600>. 1, 2, 10
- Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm’s referential dialogue magic. *ArXiv preprint*, abs/2306.15195, 2023a. URL <https://arxiv.org/abs/2306.15195>. 2
- Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, and William Cohen. MuRAG: Multimodal retrieval-augmented generator for open question answering over images and text. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5558–5570, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.375. URL <https://aclanthology.org/2022.emnlp-main.375>. 10
- Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. Can pre-trained vision and language models answer visual information-seeking questions? In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 14948–14968, Singapore, 2023b.

- Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.925. URL <https://aclanthology.org/2023.emnlp-main.925>. 1, 2, 10
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *ArXiv preprint*, abs/2404.16821, 2024. URL <https://arxiv.org/abs/2404.16821>. 2
- Wanqing Cui, Keping Bi, Jiafeng Guo, and Xueqi Cheng. MORE: Multi-mOdal REtrieval augmented generative commonsense reasoning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 1178–1192, Bangkok, Thailand and virtual meeting, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.69. URL <https://aclanthology.org/2024.findings-acl.69>. 10
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/9a6a435e75419a836fe47ab6793623e6-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/9a6a435e75419a836fe47ab6793623e6-Abstract-Conference.html). 10
- Xueying Du, Geng Zheng, Kaixin Wang, Jiayi Feng, Wentai Deng, Mingwei Liu, Bihuan Chen, Xin Peng, Tao Ma, and Yiling Lou. Vul-rag: Enhancing llm-based vulnerability detection via knowledge-level rag, 2024. URL <https://arxiv.org/abs/2406.11147>. 10
- Darryl Hannan, Akshay Jain, and Mohit Bansal. Mnymodalqa: Modality disambiguation and qa over diverse inputs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7879–7886, Apr. 2020. doi: 10.1609/aaai.v34i05.6294. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6294>. 2
- Wenbo Hu, Zi-Yi Dou, Liunian Harold Li, Amita Kamath, Nanyun Peng, and Kai-Wei Chang. Mqt-llava: Matryoshka query transformer for large vision-language models. In *The 38th Conference on Neural Information Processing Systems (NeurIPS)*, 2024a. 10
- Wenbo Hu, Yifan Xu, Yi Li, Weiyue Li, Zeyuan Chen, and Zhuowen Tu. BLIVA: A simple multimodal LLM for better handling of text-rich visual questions. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pp. 2256–2264. AAAI Press, 2024b. doi: 10.1609/AAAI.V38I3.27999. URL <https://doi.org/10.1609/aaai.v38i3.27999>. 2
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, Qiang Liu, Kriti Aggarwal, Zewen Chi, Nils Johan Bertil Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. Language is not all you need: Aligning perception with language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/e425b75bac5742a008d643826428787c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/e425b75bac5742a008d643826428787c-Abstract-Conference.html). 2
- Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. Mantis: Interleaved multi-image instruction tuning. *ArXiv preprint*, abs/2405.01483, 2024a. URL <https://arxiv.org/abs/2405.01483>. 4, 7, 10
- Ting Jiang, Minghui Song, Zihan Zhang, Haizhen Huang, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, and Fuzhen Zhuang. E5-v: Universal embeddings with multimodal large language models, 2024b. URL <https://arxiv.org/abs/2407.12580>. 9

- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, 2013. 2, 5
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. OBELICS: an open web-scale filtered dataset of interleaved image-text documents. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/e2cfb719f58585f779d0a4f9f07bd618-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/e2cfb719f58585f779d0a4f9f07bd618-Abstract-Datasets_and_Benchmarks.html). 10
- Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models?, 2024. 7
- Paul Lerner, Olivier Ferret, Camille Guinaudeau, Hervé Le Borgne, Romaric Besançon, Jose G. Moreno, and Jesus Lovon. ViQuAE, a Dataset for Knowledge-based Visual Question Answering about Named Entities. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, Madrid, Spain, July 2022. doi: 10.1145/3477495.3531753. URL <https://universite-paris-saclay.hal.science/hal-03650618>. 2
- 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. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>. 10
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024a. URL <https://arxiv.org/abs/2408.03326>. 7, 10
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench-2: Benchmarking multimodal large language models. *ArXiv preprint*, abs/2311.17092, 2023. URL <https://arxiv.org/abs/2311.17092>. 10
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models, 2024b. URL <https://arxiv.org/abs/2407.07895>. 4, 7, 10
- Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023. 7, 10
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023a. URL <https://arxiv.org/abs/2310.03744>. 2, 23
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/6dcf277ea32ce3288914faf369fe6de0-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/6dcf277ea32ce3288914faf369fe6de0-Abstract-Conference.html). 10
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *ArXiv preprint*, abs/2307.06281, 2023c. URL <https://arxiv.org/abs/2307.06281>. 10

- 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. Deepseek-vl: Towards real-world vision-language understanding, 2024a. [7](#)
- Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-GPS: Interpretable geometry problem solving with formal language and symbolic reasoning. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 6774–6786, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.528. URL <https://aclanthology.org/2021.acl-long.528>. [10](#)
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022a. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/11332b6b6c6f4485b84afadb1352d3a9a-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/11332b6b6c6f4485b84afadb1352d3a9a-Abstract-Conference.html). [10](#)
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024b. [7](#), [10](#), [23](#)
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8086–8098, Dublin, Ireland, 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.556. URL <https://aclanthology.org/2022.acl-long.556>. [9](#)
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 3195–3204. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00331. URL [http://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Marino\\_OK-VQA\\_A\\_Visual\\_Question\\_Answering\\_Benchmark\\_Requiring\\_External\\_Knowledge\\_CVPR\\_2019\\_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Marino_OK-VQA_A_Visual_Question_Answering_Benchmark_Requiring_External_Knowledge_CVPR_2019_paper.html). [2](#), [10](#)
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights from multimodal llm pre-training. *ArXiv preprint*, abs/2403.09611, 2024. URL <https://arxiv.org/abs/2403.09611>. [2](#)
- Thomas Mensink, Jasper R. R. Uijlings, Lluís Castrejón, Arushi Goel, Felipe Cadar, Howard Zhou, Fei Sha, André Araújo, and Vittorio Ferrari. Encyclopedic VQA: visual questions about detailed properties of fine-grained categories. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pp. 3090–3101. IEEE, 2023. doi: 10.1109/ICCV51070.2023.00289. URL <https://doi.org/10.1109/ICCV51070.2023.00289>. [1](#), [2](#), [10](#)
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, 2008. [2](#), [5](#)
- OpenAI. Gpt-4 technical report, 2023. [2](#), [4](#), [7](#)
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina

- Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 2021. URL <http://proceedings.mlr.press/v139/radford21a.html>. 9
- Vikram V. Ramaswamy, Sing Yu Lin, Dora Zhao, Aaron Adcock, Laurens van der Maaten, Deepti Ghadiyaram, and Olga Russakovsky. Geode: a geographically diverse evaluation dataset for object recognition. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/d08b6801f24dda81199079a3371d77f9-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/d08b6801f24dda81199079a3371d77f9-Abstract-Datasets_and_Benchmarks.html). 5
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y. 2, 5
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *European Conference on Computer Vision*, pp. 146–162. Springer, 2022. 2, 10
- Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. Kvqa: Knowledge-aware visual question answering. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):8876–8884, Jul. 2019. doi: 10.1609/aaai.v33i01.33018876. URL <https://ojs.aaai.org/index.php/AAAI/article/view/4915>. 2
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. REPLUG: Retrieval-augmented black-box language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8371–8384, Mexico City, Mexico, 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.naacl-long.463>. 10
- Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, et al. Generative multimodal models are in-context learners. *ArXiv preprint*, abs/2312.13286, 2023. URL <https://arxiv.org/abs/2312.13286>. 10
- Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. Multimodal{qa}: complex question answering over text, tables and images. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=ee6W5UgQLa>. 2
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *ArXiv preprint*, abs/2312.11805, 2023. URL <https://arxiv.org/abs/2312.11805>. 4, 7
- Mistral AI Team. Announcing pixtral 12b. <https://mistral.ai/news/pixtral-12b/>, 2024. 7
- Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms, 2024. 2
- Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language models are latent variable models: Explaining and finding good

- demonstrations for in-context learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/3255a7554605a88800f4e120b3a929e1-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/3255a7554605a88800f4e120b3a929e1-Abstract-Conference.html). 9
- Keyang Xuan, Li Yi, Fan Yang, Ruochen Wu, Yi R. Fung, and Heng Ji. Lemma: Towards lvlm-enhanced multimodal misinformation detection with external knowledge augmentation, 2024. URL <https://arxiv.org/abs/2402.11943>. 10
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Richard James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-Tau Yih. Retrieval-augmented multimodal language modeling. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 39755–39769. PMLR, 2023. URL <https://proceedings.mlr.press/v202/yasunaga23a.html>. 10
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models, 2024. URL <https://arxiv.org/abs/2408.04840>. 7
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *ArXiv preprint*, abs/2306.13549, 2023. URL <https://arxiv.org/abs/2306.13549>. 10
- Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *ArXiv preprint*, abs/2404.16006, 2024. URL <https://arxiv.org/abs/2404.16006>. 10
- Zheng Yuan, Qiao Jin, Chuanqi Tan, Zhengyun Zhao, Hongyi Yuan, Fei Huang, and Songfang Huang. Ramm: Retrieval-augmented biomedical visual question answering with multi-modal pre-training. *Proceedings of the 31st ACM International Conference on Multimedia*, 2023. 10
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *ArXiv preprint*, abs/2311.16502, 2023. URL <https://arxiv.org/abs/2311.16502>. 10
- Kai Zhang, Yi Luan, Hexiang Hu, Kenton Lee, Siyuan Qiao, Wenhui Chen, Yu Su, and Ming-Wei Chang. Magiclens: Self-supervised image retrieval with open-ended instructions. In *The Forty-first International Conference on Machine Learning (ICML)*, pp. to appear, 2024a. 9
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? *ArXiv preprint*, abs/2403.14624, 2024b. URL <https://arxiv.org/abs/2403.14624>. 10
- Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruva Ghosh, Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy. Why are visually-grounded language models bad at image classification? *ArXiv preprint*, abs/2405.18415, 2024c. URL <https://arxiv.org/abs/2405.18415>. 5
- Ruochen Zhao, Hailin Chen, Weishi Wang, Fangkai Jiao, Xuan Long Do, Chengwei Qin, Bosheng Ding, Xiaobao Guo, Minzhi Li, Xingxuan Li, and Shafiq Joty. Retrieving multimodal information for augmented generation: A survey. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 4736–4756, Singapore, 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.314. URL <https://aclanthology.org/2023.findings-emnlp.314>. 10



- Junjie Zhou, Zheng Liu, Shitao Xiao, Bo Zhao, and Yongping Xiong. Vista: Visualized text embedding for universal multi-modal retrieval, 2024. URL <https://arxiv.org/abs/2406.04292>. 9
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *ArXiv preprint, abs/2304.10592*, 2023a. URL <https://arxiv.org/abs/2304.10592>. 10
- Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal C4: an open, billion-scale corpus of images interleaved with text. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/1c6bed78d3813886d3d72595dbecb80b-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/1c6bed78d3813886d3d72595dbecb80b-Abstract-Datasets_and_Benchmarks.html). 10

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>MRAG-BENCH</b>	<b>3</b>
2.1	Benchmark Overview . . . . .	3
2.2	Benchmark Composition . . . . .	4
2.3	Data Collection . . . . .	5
<b>3</b>	<b>Experiments</b>	<b>7</b>
3.1	Experimental Setup . . . . .	7
3.2	Main Results . . . . .	7
<b>4</b>	<b>Analysis</b>	<b>8</b>
4.1	How much can visual knowledge benefit more than textual knowledge? . . . . .	8
4.2	How does retriever performance affect LVLMs? . . . . .	9
4.3	How many ground-truth image examples are needed? . . . . .	9
<b>5</b>	<b>Related Work</b>	<b>10</b>
5.1	Multimodal Retrieval-Augmented Generation Benchmarks . . . . .	10
5.2	Large Vision Language Models . . . . .	10
5.3	Retrieval-Augmented Multimodal Large Language Models . . . . .	10
<b>6</b>	<b>Conclusion</b>	<b>10</b>
<b>A</b>	<b>MRAG-BENCH Details</b>	<b>18</b>
A.1	Dataset Curation Details . . . . .	18
A.2	Human Evaluation Protocol . . . . .	22
<b>B</b>	<b>Experiment Setting Details</b>	<b>23</b>
B.1	Model Prompts . . . . .	23
B.2	Evaluation Tool . . . . .	23
<b>C</b>	<b>More Results</b>	<b>24</b>

## A MRAG-BENCH DETAILS

## A.1 DATASET CURATION DETAILS

**Dataset collection of transformative aspect** We chose to manually scrape images from the web based on the definitions of the transformative aspect. To construct the image corpus, we employed Bing Image Search for each of the image object keyword predefined by us. We filtered some of the search results where the image objects do not have a clear pair of query image and ground-truth image example, around 74% keyword names were kept during this process. Here we listed all the keywords that are already filtered and used for search of query image except in biological scenario, it’s for search of ground-truth image example. Each search keyword is composed of an “image

object” and a “condition”. For example, “A young kitten image of Himalayan Cat”, here Himalayan Cat is the image object and a young kitten is the condition. For each of keyword listed below, we searched again for its ground-truth examples (except for biological scenario, it’s for query images), in which only “image object” is kept and “condition” is removed. All searched results are further picked and downloaded by humans to ensure quality. Here is a list of the filtered keywords for transformative aspect:

**Transformative: Temporal**

- A young kitten image of Himalayan Cat
- A young kitten image of Chartreux
- A young kitten image of Burmese
- A young kitten image of Turkish Van
- A young kitten image of American Shorthair
- A young kitten image of British Shorthair
- A young kitten image of Maine Coon
- A young kitten image of Burma (Myanmar)
- A young kitten image of Selkirk Rex
- A young kitten image of Siberian
- A young kitten image of Persian
- A young kitten image of Manx
- A young kitten image of Ocicat
- A young kitten image of Russian Blue
- A young kitten image of Bengal Cat
- A young kitten image of Devon Rex
- A young kitten image of American Bobtail
- A young kitten image of Balinese
- A young kitten image of LaPerm
- A young kitten image of Egyptian Mau
- A young kitten image of Japanese Bobtail
- A young kitten image of Ragdoll
- A young kitten image of Abyssinian
- A young kitten image of American Wirehair
- A young kitten image of Oriental Shorthair
- A young kitten image of Cornish Rex
- A young kitten image of Kurilian Bobtail
- A young kitten image of Singapura Cat
- A young kitten image of Birman
- A young kitten image of Burmilla
- A young kitten image of Korat
- A young kitten image of Tonkinese
- A young kitten image of Somali Cat
- A young kitten image of Norwegian Forest Cat
- A young kitten image of Turkish Angora
- A young kitten image of Siamese
- A picture of Sainte-Chapelle under construction
- A picture of Washington Monument under construction
- A picture of Hearst Castle under construction
- A picture of Time Square under construction
- A picture of Wrigley Building under construction
- A picture of Eiffel Tower under construction
- A picture of The Arc de Triomphe under construction
- A picture of Golden Gate Bridge under construction
- A picture of White House under construction
- A picture of Palace of Versailles under construction
- A picture of Opéra Garnier under construction
- A picture of San Simeon under construction
- A picture of The Louvre under construction
- A picture of Cathédrale Notre-Dame de Paris under construction
- A picture of Sacré-Cœur Basilica under construction

- A picture of Brooklyn Bridge under construction
- A picture of Panthéon under construction
- A picture of Capitol Building under construction
- A picture of Independence Hall under construction
- A picture of Mont Saint-Michel under construction
- A picture of St Patrick's Cathedral under construction
- A picture of Space Needle under construction
- A picture of Château de Chambord under construction
- A picture of Versailles under construction

**Transformative: Deformation**

- An image of Toyota Camry damaged
- An image of Ford F-150 damaged
- An image of Ferrari 458 damaged
- An image of Audi Q5 damaged
- An image of Lamborghini LP640 damaged
- An image of McLaren 675LT damaged
- An image of Mercedes SLC damaged
- An image of Lamborghini Aventador damaged
- An image of Lamborghini LP570 damaged
- An image of Porsche 911 GT3 RS damaged
- An image of Audi A6 damaged
- An image of Audi A4 damaged
- An image of Lamborghini Aventador SV damaged
- An image of GMC Sierra 2500 HD damaged
- An image of Infiniti G37 damaged
- An image of GMC Yukon damaged
- An image of Honda Accord damaged
- An image of Infiniti FX35 damaged
- An image of Tesla Model 3 damaged
- An image of Acura RDX 2020 damaged
- An image of BMW 7 Series damaged
- An image of Audi A5 Sportback damaged
- An image of Hyundai IX35 damaged
- An image of Cadillac XTS damaged
- An image of BMW M3 damaged
- An image of Acura MDX damaged
- An image of Audi A3 damaged
- An image of BMW X3 damaged
- An image of Porsche Boxster damaged
- An image of Mercedes CLA45 AMG damaged
- An image of Jaguar XJ damaged

**Transformative: Incomplete**

- MacBook Keyboard missing keys
- Windows Keyboard missing keys
- Laptop Keyboards (Generic) missing keys
- Mechanical Keyboard missing keys
- Ergonomic Keyboard missing keys
- Compact Keyboard missing keys
- Gaming Keyboard missing keys
- Chiclet Keyboard missing keys
- Tenkeyless (TKL) Keyboard missing keys
- Virtual Keyboard (On-screen) missing keys
- Numeric Keypad missing keys
- ISO Keyboard Layout missing keys
- ANSI Keyboard Layout missing keys
- Ortholinear Keyboard missing keys
- Bluetooth/Wireless Keyboard missing keys

**Transformative: Biological**

- An image of Lime after oxidation
- An image of breadfruit after oxidation
- An image of dragonfruit after oxidation
- An image of starfruit after oxidation
- An image of Raspberry after oxidation
- An image of Zucchini after oxidation
- An image of Pear after oxidation
- An image of passionfruit after oxidation
- An image of Blackberry after oxidation
- An image of durian after oxidation
- An image of persimmon after oxidation
- An image of Apple after oxidation
- An image of bell pepper after oxidation
- An image of olive after oxidation
- An image of Mango after oxidation
- An image of nectarine after oxidation
- An image of tomato after oxidation
- An image of quince after oxidation
- An image of coconut after oxidation
- An image of soursop after oxidation
- An image of Kiwi after oxidation
- An image of cucumber after oxidation
- An image of apricot after oxidation
- An image of Honeydew after oxidation
- An image of Peach after oxidation
- An image of pomegranate after oxidation
- An image of carrot after oxidation
- An image of fig after oxidation
- An image of Papaya after oxidation
- An image of Blueberry after oxidation
- An image of Banana after oxidation
- An image of jackfruit after oxidation
- An image of Lemon after oxidation
- An image of tamarind after oxidation
- An image of lychee after oxidation
- An image of Pineapple after oxidation
- An image of Cantaloupe after oxidation
- An image of Orange after oxidation
- An image of Rambutan after oxidation
- An image of guava after oxidation
- An image of sweet potato after oxidation
- An image of Plum after oxidation
- An image of Avocado after oxidation
- An image of Watermelon after oxidation
- An image of potato after oxidation
- An image of Grapefruit after oxidation
- An image of Grapes after oxidation
- An image of pumpkin after oxidation
- An image of Cherry after oxidation
- An image of Strawberry after oxidation
- An image of custard apple after oxidation

**Quality control** We employ two types of quality control throughout the annotation process: an automatic check with predefined rules and a manual examination of each instance. The automatic check verifies correct MCQA format in which each question should only have one correct answer, metadata values, assesses image validity (checking the accessibility of each image) and filters out redundant images in the corpus (images that are repetitively downloaded). The manual examination

is conducted by two experts in this field, who checked the correspondence between query images and ground-truth image examples, and filtered or revised ambiguous questions and uncorrelated query image and ground-truth images.

## A.2 HUMAN EVALUATION PROTOCOL

Three human annotators in domain conducted the human evaluation. The interface for human evaluation without RAG knowledge and with RAG knowledge are shown in Figure 6 and Figure 7.

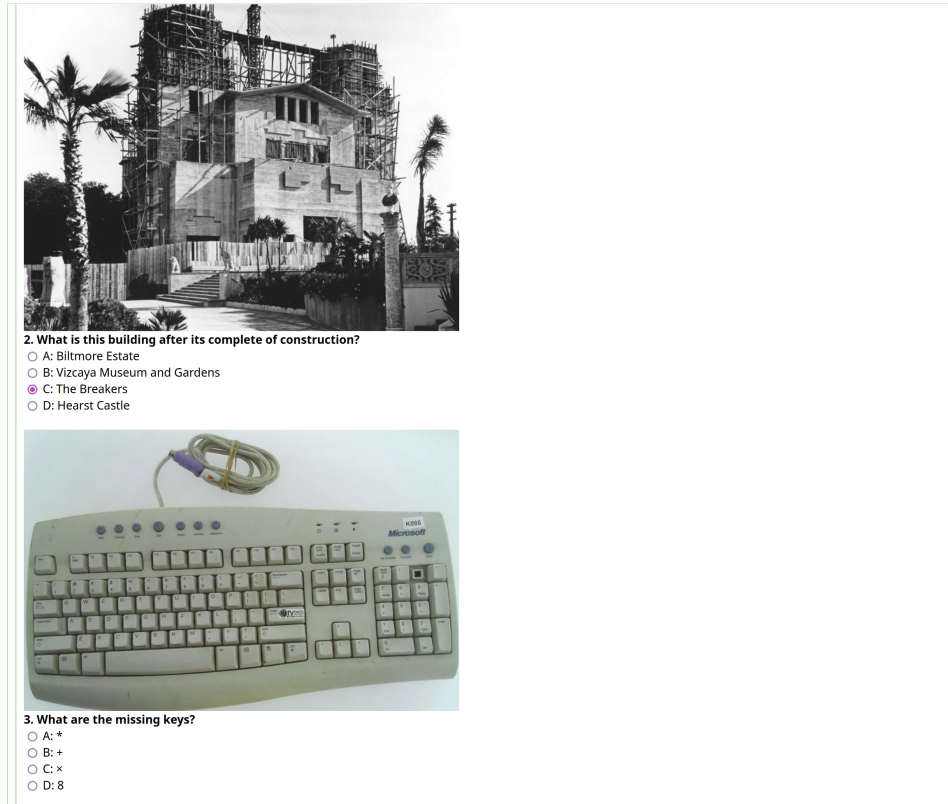


Figure 6: Human evaluation interface without RAG examples

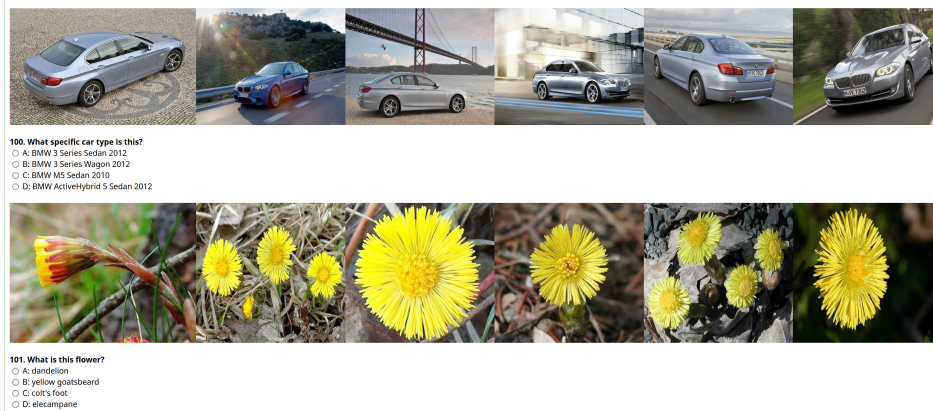


Figure 7: Human evaluation interface with ground-truth RAG examples

## B EXPERIMENT SETTING DETAILS

### B.1 MODEL PROMPTS

Following Lu et al. (2024b) and Liu et al. (2023a) our prompt consists of four parts, the instruction, question, options, and a prefix of the answer. For images, we insert them into the text to form a coherent prompt as the image placeholder (`{Image}`) indicated below. The complete prompt is as follows:

Model Prompts for No RAG Evaluation
Instruction: Answer with the option’s letter from the given choices directly. <code>{Image}</code> Question: <code>{QUESTION}</code> Choices: (A) <code>{OPTION_A}</code> (B) <code>{OPTION_B}</code> (C) <code>{OPTION_C}</code> (D) <code>{OPTION_D}</code> Answer:
Model Prompts for RAG Evaluation
Instruction: You will be given one question concerning several images. The first image is the input image, others are retrieved examples to help you. Answer with the option’s letter from the given choices directly. <code>{Image}{Image}{Image}{Image}{Image}{Image}</code> Question: <code>{QUESTION}</code> Choices: (A) <code>{OPTION_A}</code> (B) <code>{OPTION_B}</code> (C) <code>{OPTION_C}</code> (D) <code>{OPTION_D}</code> Answer:

### B.2 EVALUATION TOOL

Following Lu et al. (2024b), we first use a rule-based automatic tool to extract the exact answer. First, the tool detects if a valid option index appears in the model output. If no direct answer is found, the tool matches the output to the content of each option. If there is still no match, we employ GPT-3.5-turbo to automatically extract the answer following our prompts in Table 5. If GPT-3.5-turbo finds there is still no match, we will randomly select an option as the answer.

<p><b>Prompt</b>            Please read the following example. Then extract the multiple choice letter with the answer corresponding to the choice list from the model response and type it at the end of the prompt. You should only output either A, B, C, or D.</p> <p><code>{In-context examples}</code></p> <p>Question: <code>{QUESTION}</code>            Choice List: (A) <code>{OPTION_A}</code> (B) <code>{OPTION_B}</code> (C) <code>{OPTION_C}</code> (D) <code>{OPTION_D}</code>            Model Response: <code>{Response}</code>            Extracted answer:</p>
---

Table 5: Prompt template to extract multiple choice answer from model’s response. `{In-context examples}` are in-context examples.

## C MORE RESULTS

We present the Recall@5 scores per each scenarios on 4 multimodal retrievers as shown in Table 6 and LLaVA-Next-Interleave’s accuracy score affected by these retrievers in Table 7.

Model	Overall	Perspective				Transformative				Others
		Angle	Partial	Scope	Occlusion	Temporal	Deformation	Incomplete	Biological	
MagicLens	37.03	41.61	33.33	36.27	36.11	12.75	10.78	79.41	29.41	56.67
E5-V	54.92	49.69	48.78	61.76	66.67	38.93	22.55	73.53	71.57	82.50
VISTA	59.65	66.15	67.48	64.71	63.89	38.26	8.82	33.33	94.12	80.83
CLIP	60.46	70.19	54.47	71.57	73.15	44.30	31.37	67.65	40.2	81.67

Table 6: Recall@5 scores with 4 retriever models on MRAG-BENCH.

Model	Overall	Perspective				Transformative				Others
		Angle	Partial	Scope	Occlusion	Temporal	Deformation	Incomplete	Biological	
MagicLens	35.18	34.78	29.67	30.39	34.26	40.94	36.27	27.45	49.02	39.17
E5-V	40.06	38.82	39.84	41.18	46.3	38.93	41.18	27.45	48.04	41.67
VISTA	42.42	40.37	35.77	40.2	52.78	45.64	42.16	36.27	50.98	48.33
CLIP	40.35	40.06	33.33	39.22	56.48	43.62	44.12	27.45	36.27	49.17

Table 7: LLaVA-Next-Interleave accuracy scores on MRAG-BENCH with 4 different retrievers.