



# MMSEARCH: BENCHMARKING THE POTENTIAL OF LARGE MODELS AS MULTI-MODAL SEARCH ENGINES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The advent of Large Language Models (LLMs) has paved the way for AI search engines, e.g., SearchGPT, showcasing a new paradigm in human-internet interaction. However, most current AI search engines are limited to text-only settings, neglecting the multimodal user queries and the text-image interleaved nature of website information. Recently, Large Multimodal Models (LMMs) have made impressive strides. Yet, whether they can function as AI search engines remains under-explored, leaving the potential of LMMs in multimodal search an open question. To this end, we first design a delicate pipeline, **MMSEARCH-ENGINE**, to empower any LMMs with multimodal search capabilities. On top of this, we introduce **MMSEARCH**, a comprehensive evaluation benchmark to assess the multimodal search performance of LMMs. The curated dataset contains 300 manually collected instances spanning 14 subfields, which involves no overlap with the current LMMs’ training data, ensuring the correct answer can only be obtained within searching. By using **MMSEARCH-ENGINE**, the LMMs are evaluated by performing three individual tasks (requery, rerank, and summarization), and one challenging end-to-end task with a complete searching process. We conduct extensive experiments on closed-source and open-source LMMs. Among all tested models, GPT-4o with **MMSEARCH-ENGINE** achieves the best results, which surpasses the commercial product, Perplexity Pro, in the end-to-end task, demonstrating the effectiveness of our proposed pipeline. We further present error analysis to unveil current LMMs still struggle to fully grasp the multimodal search tasks, and conduct ablation study to indicate the potential of scaling test-time computation for AI search engine. We hope **MMSEARCH** may provide unique insights to guide the future development of multimodal AI search engine.

## 1 INTRODUCTION

Search engines (Brin & Page, 1998) have been the main tools for humans to navigate through the overwhelming quantity of online resources. Recently, Large Language Models (LLMs) (OpenAI, 2023a;b; Touvron et al., 2023a) have demonstrated impressive performance on various zero-shot downstream applications. On top of this, AI search engine (OpenAI, 2024c), which integrates LLMs with traditional search engines, stands among one of the most promising ones. It points the direction of the next-generation interaction paradigm of human and Internet. Combining the language understanding ability of LLMs and up-to-date information from the Internet, AI search engines could better grasp the user’s intention and summarize contextual-aligned answers from the raw web information. These systems can only process textual queries and interpret textual web content, significantly constraining user query scenarios and information-seeking methods (Barbany et al., 2024; Xie et al., 2024). This limitation impacts both the range of input queries and the accuracy of results (Jiang et al., 2024a; Chen et al., 2021; Lù et al., 2024), particularly given the complexity and interleaved nature of modern websites (Liu et al., 2024b). For example, consider a scenario where you possess numerous medals belonging to your grandfather but are unaware of their specific names. A multimodal AI search engine could match photographs of these medals with an interleaved table of images and text retrieved from the Internet, thereby identifying each medal. In contrast, text-

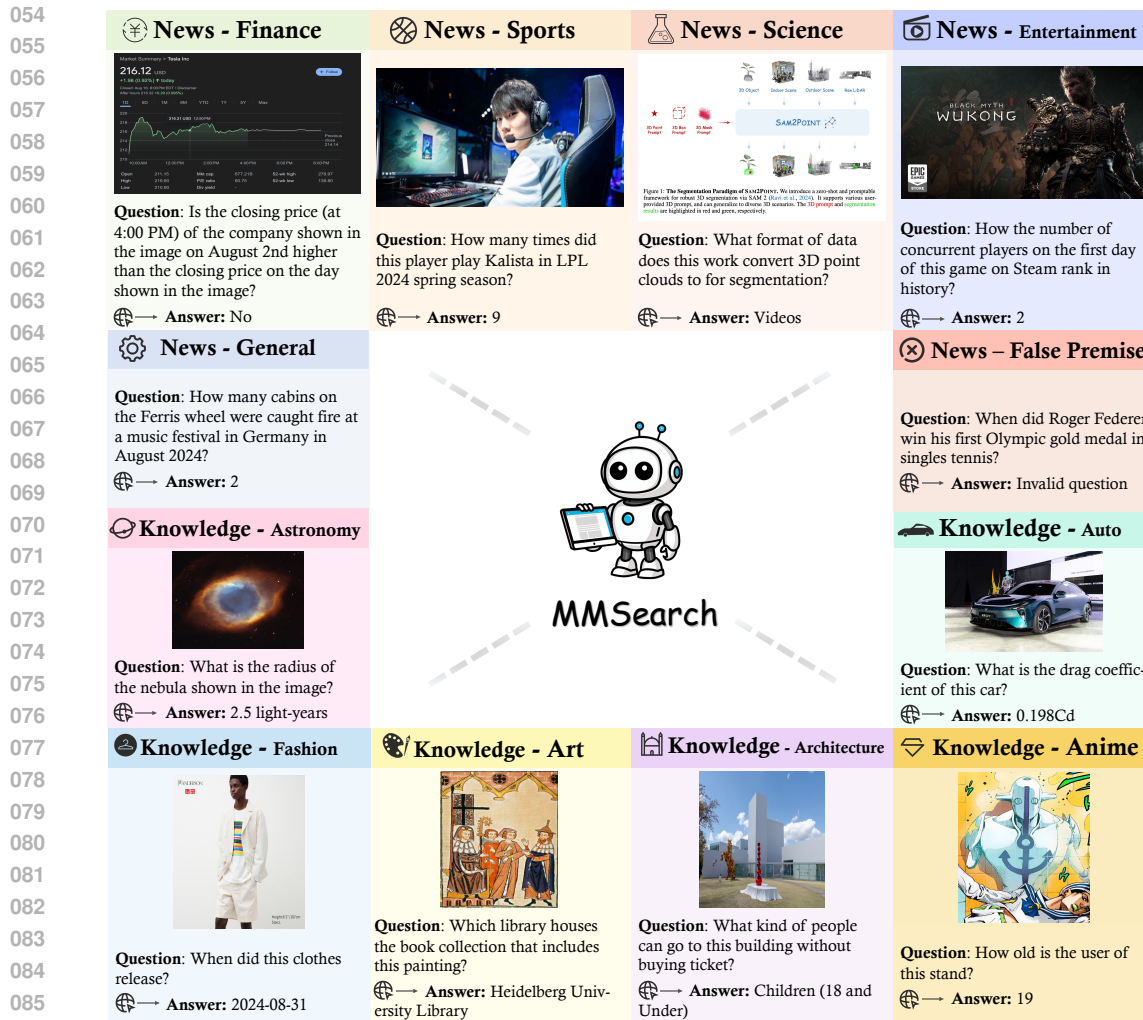


Figure 1: **Overview of the MMSEARCH Benchmark.** MMSEARCH aims to evaluate any LMM’s potential to be a multimodal AI search engine. The benchmark contains two primary areas: latest news and rare knowledge to ensure no overlap with LMM’s inherent knowledge.

only search engines can neither take photographs for searching nor understand the interleaved table. Hence, a multimodal AI search engine is crucial for advancing information retrieval and analysis.

On the other hand, with the recent rapid advancements, Large Multimodal Models (LMMs) (Liu et al., 2023a; Lin et al., 2023; OpenAI, 2023c; Gao et al., 2024; Zhang et al., 2024b) have showcased significant abilities across diverse scenarios, including general image understanding (Fu et al., 2023; Liu et al., 2023b; Yu et al., 2023), expert image reasoning (Zhang et al., 2024d; Gao et al., 2023a; Zhang et al., 2024c; Guo et al., 2024a), multi-image perception (Li et al., 2024a; Wang et al., 2024; Jiang et al., 2024b; Li et al., 2024c), and spatial environment perception (Guo et al., 2023; Yang et al., 2023; Han et al., 2023). Despite these developments, a framework for LMMs to function as multimodal AI search engines remains largely unexplored. Consequently, the potential of LMMs in multimodal searching also remains a significant open question.

To bridge this gap, we first present MMSEARCH-ENGINE, a multimodal AI search engine pipeline, empowering any LMMs with advanced search capabilities. MMSEARCH-ENGINE maximizes the utilization of LMMs’ multimodal information comprehension abilities, incorporating both visual and textual website content as information sources. On top of this, we introduce MMSEARCH, a multimodal AI search engine benchmark to comprehensively evaluate LMMs’ searching performance. The design of MMSEARCH-ENGINE facilitates the zero-shot evaluation of any LMMs

108 within the context of AI search engine. Our experiment covers state-of-the-art closed-source (OpenAI, 2023c; Anthropic, 2024; Gemini Team, 2023) and open-source LMMs (Li et al., 2024b; Qwen Team, 2024; Chen et al., 2024d; Ye et al., 2024). Our efforts are summarized as follows:

- 112 i. **MMSEARCH-ENGINE, a multimodal AI search engine pipeline for LMMs**, empowering  
113 large models for multimodal searching. In contrast with the conventional text-only  
114 AI search engines, MMSEARCH-ENGINE fully integrates multimodal information in two  
115 ways: (i) for queries containing images, we conduct web searches across both textual and  
116 visual modalities. We utilize Google Lens ([len](#)) to identify critical visual information from  
117 the input image; (ii) all search results are presented in both textual and visual formats,  
118 ensuring a comprehensive understanding of the interleaved website content. The working  
119 flow of MMSEARCH-ENGINE contains multi-round interaction between LMM and the In-  
120 ternet. The LMM needs to first *requery* the user question into a search-engine-friendly  
121 format. Then, the LMM *reranks* the retrieved websites based on its helpfulness. Finally,  
122 the LMM is required to *summarize* the answer based on the most informative webpage con-  
123 tent selected from the rerank. Thanks to the design of the pipeline, we propose a step-wise  
124 evaluation strategy on the three core tasks within the searching process: *requery*, *rerank*,  
125 and *summarization*. The final score is weighted by the end-to-end evaluation results and  
126 scores of the three core tasks.
- 127 ii. **MMSEARCH, a comprehensive benchmark for multimodal AI search engines**, which, to  
128 our best knowledge, serves as the first evaluation dataset to measure LMMs’ multimodal  
129 searching capabilities. Our benchmark categorizes searching queries into two primary ar-  
130 eas: *News* and *Knowledge*, as shown in Fig. 1. We employed different strategies for these  
131 two areas to ensure the challenging nature of the benchmark. *News* area covers the latest  
132 news at the time of data collection (August, 2024). This is to guarantee the answers to  
133 the queries will not be present in the training data of LMMs. As for the area of *Knowl-  
134 edge*, we collect queries requiring rare knowledge and then select the queries unable to be  
135 answered by current SoTA LMMs such as GPT-4o (OpenAI, 2024b) or Claude-3.5 (An-  
136 thropic, 2024). The two areas sum up to 14 subfields. In total, MMSEARCH encompasses  
137 300 meticulously collected queries, with 2901 unique images.
- 138 iii. **Extensive experiments and error analysis for future development direction**. We evaluate  
139 popular closed-source models and open-source LMMs on MMSEARCH. GPT-4o achieves  
140 the best overall performance across different tasks. Surprisingly, our MMSEARCH-  
141 ENGINE equipped with SoTA LMMs, such as GPT-4o and Claude 3.5 Sonnet, even sur-  
142 passes the prominent commercial AI search engine Perplexity Pro (Perplexity) in the end-  
143 to-end task. Our thorough error analysis reveals that current LMMs still struggle to gen-  
144 eralize to multimodal search-specific tasks. Their poor requery and rerank capabilities  
145 significantly limit their ability to correctly identify useful websites and extract relevant  
146 answers. Additionally, we identify five error types for requery and summarization tasks,  
147 respectively. We find that current LMMs cannot fully understand the requery task and do  
148 not know how to query the search engine. As for the summarization task, LMMs often  
149 have difficulty in extracting useful information, either from text or images. These capabil-  
150 ities are essential for LMMs to function as robust multimodal search engines and require  
151 further development. We also conduct a preliminary ablation study to explore the potential  
152 of scaling test-time computation versus scaling model size (OpenAI, 2024a). Initial results  
153 indicate that scaling test-time computation demonstrates superior performance in this task.

## 154 2 MMSEARCH

158 In Section 2.1, we first detail the design of our multimodal AI search engine pipeline, which serves  
159 as both data collection and evaluation tools. Then, in Section 2.2, we detail the data composition  
160 and collection of the curated multimodal search benchmark MMSEARCH. Then, in Section 2.3,  
161 we elaborate on our step-wise evaluation strategy. Finally, we detail the dynamic nature of our  
benchmark in Section 2.4.

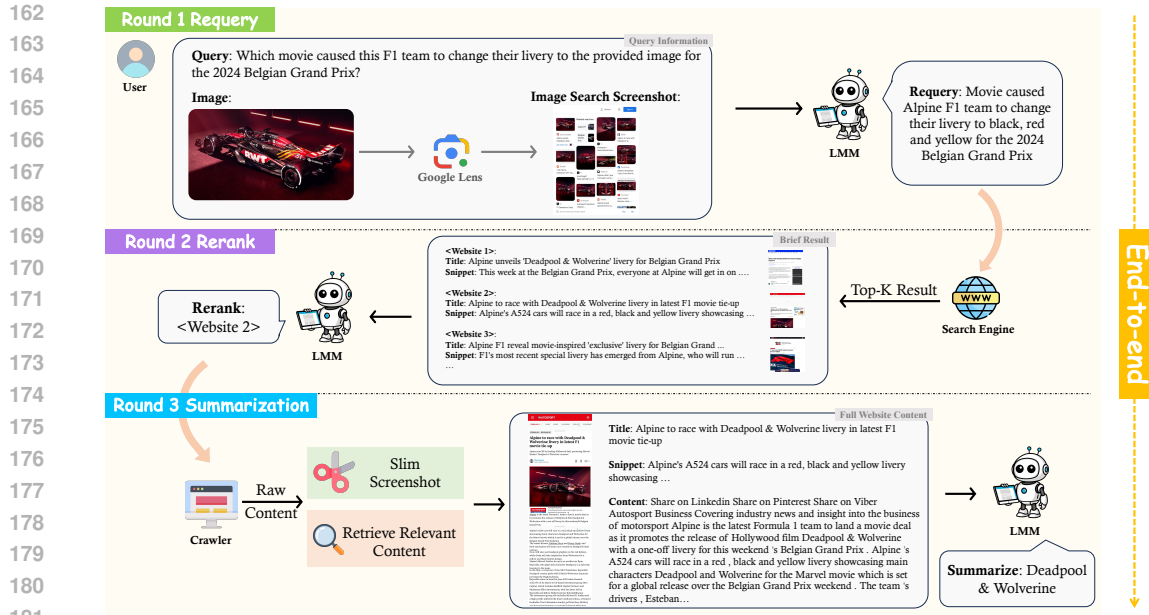


Figure 2: **The Pipeline of MMSEARCH-ENGINE.** The process comprises three sequential stages executed by a LMM: (i) query, (ii) rerank, and (iii) summarization. In the end-to-end evaluation task, the LMM completes these three stages sequentially to generate the final output.

## 2.1 MMSEARCH-ENGINE: A MULTIMODAL AI SEARCH ENGINE PIPELINE

The searching process is a complex action including multi-round interactions between LMMs and conventional search engines. We develop a delicate pipeline that queries LMMs multiple times to accomplish this task. Leveraging the image comprehension capabilities of LMMs, we incorporate two types of visual data. First, we incorporate Google Lens (*len*) to search for information from the image. The second type of visual data is the screenshot of the retrieved websites, in the purpose of preserving the original format of website content. Our framework is shown in Fig. 2. Below we detail how an LMM works with this pipeline, which comprises three sequential phases:

- i. *Requry*. The query direct from users may contain references to certain information in the image, e.g., the *News-Finance* example shown in Fig. 1. Since a conventional search engine only accept text-only input, it is necessary for LMM to translate the image content and combine it with the query to ask a valid question to it. In addition, the raw user query may be ambiguous or inefficient sometimes (Chan et al., 2024; Ma et al., 2023), reformulating the query to be more clear is also a must for LMM. If the user query contains an image, we incorporate the screenshot of the image search result from the google lens (*len*). We treat the user query, user image, and the image search screenshot as basic information of the query. This information will be input to LMM in every round in the pipeline. For the requry round, we prompt LMM to output a requry to a conventional search engine.
- ii. *Rerank*. The requry is sent to a search engine API, e.g., DuckDuckGo, to retrieve top  $K$  relevant websites. Depending on the requry quality, not all retrieved websites are necessarily relevant for query answering. Hence, we prompt LMM to select one most informative website for answer summarization. Due to the LMM’s context length limitations and the extensive content of websites, we provide only essential information of each website, which we term *brief results*. These brief results include the title, the snippet, and a screenshot of the webpage’s top section, which serves as the input for LMM’s reranking. The inclusion of the screenshot serves two purposes. First, the screenshot offers a visual cue to assess the web’s credibility, as a well-organized website often appears more trustworthy than one cluttered with advertisements (Fogg et al., 2001; Sillence et al., 2004). Additionally, the screenshot may contain essential visual information. For instance, it might include images similar or identical to query images, as shown in the *Website 2* in Fig. 2.

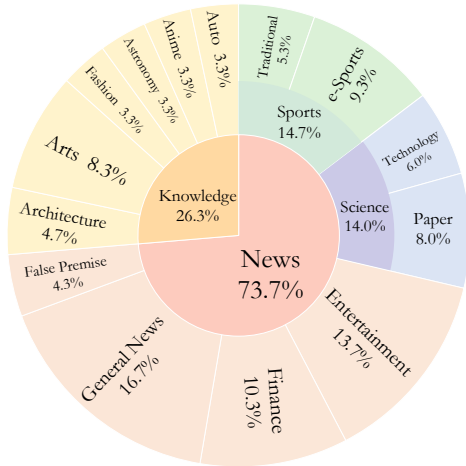


216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235

Table 1: Key Statistics of MMSEARCH.

Statistic	Number
Total questions	300
- Questions with images	171 (57.0%)
- Questions without images	129 (43.0%)
Total Websites	2,280
Total Areas/Subfields	2/14
Number of unique images	2,901
- Query images	163
- Google search images	163
- Top section screenshot images	2,280
- Full-page screenshot images	295
Number of unique questions	289
Number of unique requeries	289
Number of unique reranked websites	2,400
Number of unique answers	264
Maximum question length	41
Maximum answer length	12
Average question length	14.0
Average answer length	1.9

Figure 3: Area and Subfield Distribution of MMSEARCH.



236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249

iii. *Summarization.* We start by crawling the selected website to gather all the available information. We parse the HTML to obtain the raw textual content and capture a full-page screenshot of the website. However, there are two issues: the raw content tends to be extensively lengthy and disorganized, while substantial areas in the full-page screenshot are blank due to the ad blocks on the website. These two issues lead to a large number of input tokens filled with irrelevant information. To enhance data efficiency, we slim the screenshot and retrieve the relevant content before inputting them to LMM. For the full-page screenshot, we identify the blank areas and remove them iteratively, detailed in Appendix F. As for the text content, we apply a text embedding model (Chen et al., 2024a) to retrieve a maximum of 2K tokens relevant to the requery from the raw content. We define the slimmed screenshot and the retrieved content as *full website content*. Finally, we input the full website content, website title, and website snippet, along with the query information, to LMM for summarizing the answer.

250

2.2 DATA COMPOSITION AND COLLECTION

251

To thoroughly assess multimodal search proficiency, we compile a comprehensive problem set covering a broad spectrum of news topics, specialized knowledge domains, and query image patterns. This widespread collection for MMSEARCH aims to simulate diverse user searching scenarios, ensuring a robust evaluation of LMMs’ capabilities in multimodal search.

252

**Data Composition and Categorization.** Our benchmark aims to isolate LMMs’ inherent knowledge and assess their actual search capabilities. We focus on two primary areas: *News* and *Knowledge*. For the *News* area, the queries are related to the latest news at the time of data collection (August 2024). This guarantees no overlap between the current LMMs’ training data and questions in our benchmark. All questions in this area are recorded with their occurrence time. For fairness, LMMs with recently updated knowledge should be tested on queries that occurred after their latest data update. Due to its time-sensitive nature, the *News* area serves as a dynamic part of our benchmark. Please refer to Section 2.4 for details. As for the *Knowledge* area, we focus on rare knowledge in targeted domains. Each question proposed by an annotator is verified to be beyond the capabilities of state-of-the-art Large Language Models (LLMs) such as GPT-4o (OpenAI, 2024b) or Claude 3.5 Sonnet (Anthropic, 2024). The *Knowledge* area serves as a static component of our benchmark and remains constant over time. We collect a total of 300 queries across the 2 primary areas and 14 subfields. This dataset size balances comprehensive evaluation with efficiency, considering the multi-round interactions with the search engine and Internet for each query. Detailed statistics for data composition and categorization are presented in Table 1 and Fig. 3. Definitions of each subfield are in Appendix G.2.

263  
264  
265  
266  
267  
268  
269

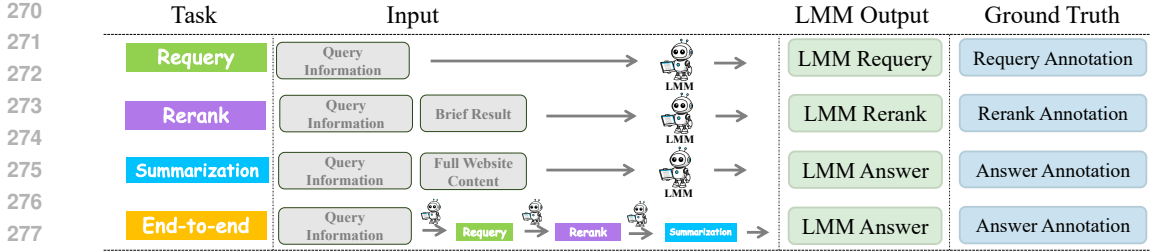


Figure 4: **Outline of Evaluation Tasks, Inputs, and Outputs.** Our evaluation contains four tasks. The requery, rerank, and summarization tasks assess the LMM’s proficiency in individual pipeline rounds. The end-to-end task simulates a real-world search scenario by sequentially executing all three stages. An example of the input and output is shown in Fig. 11 in the supplementary material.

**Data Collection and Review Process.** Thanks to the design of our pipeline, the data collection process follows similar procedure introduced in the pipeline. An annotator is first required to propose a query and provide its answer, either sourced from the latest news or rare knowledge. The annotator then formulates a requery based on the query information. After  $K$  websites are retrieved from the search engine, the annotator is required to divide all  $K$  websites into three sets based on the brief results: *valid* (likely to contain the answer), *unsure* (relevance is difficult to determine), and *invalid* (entirely irrelevant to the question). We mandate that at least one website must be classified as *valid*; if this criterion is not met, the annotator is required to adjust the requery to obtain new search results. Finally, we randomly pick one website from the *valid* set and obtain its full content. To ensure the question is answerable, another annotator is employed to give an answer to the query based on the full content. If the answer is incorrect, the question needs to be revised until it is answerable.

### 2.3 EVALUATION PROTOCOL

In contrast with previous LMM benchmarks, the multimodal search process of LMM contains multiple rounds. Only the end-to-end evaluation of the final answer is inadequate to reveal the models’ deficiency in each core searching step. For example, the errors made by the model may occur during the summarization process, but it might also stem from choosing an incorrect website during the reranking stage. To this end, we propose a step-wise strategy to evaluate the LMMs’ capability on the three core searching steps, in addition to the end-to-end evaluation.

- **End-to-end score ( $S_{e2e}$ ):** We compute the F1 score between the predicted answer and the ground truth to judge if the answer is correct.
- **Requery score ( $S_{req}$ ):** We apply the average of ROUGE-L and BLEU-1 scores to measure the similarity between the model’s requery and human-annotated requery.
- **Rerank score ( $S_{rer}$ ):** The rerank score is derived from the LMM’s selection among  $K$  pre-defined websites. The score values is 1.0 for valid set, 0.5 for unsure set, and 0 for invalid set or incorrect format.
- **Summarization score ( $S_{sum}$ ):** Again, we compute the F1 score of LMM’s answer based on a pre-defined website content against ground truth.

The input, output, and ground truth of the four tasks are visualized in Fig. 4. The final score is weighted by these four scores. We assign the highest weight (75%) to the end-to-end task, as it reflects the real-world multimodal search capability. The remaining 25% is distributed among the intermediate steps: 10% each for the rerank and summarization tasks, and 5% for the requery task. The lower weight for the requery task accounts for the inherent uncertainty in this process. The scoring process can be formulated as:

$$S_{final} = 0.75 \cdot S_{e2e} + 0.05 \cdot S_{req} + 0.1 \cdot S_{rer} + 0.1 \cdot S_{sum} \tag{1}$$

### 2.4 BENCHMARK EVOLUTION

In Fig. 5, we showcase the statistics of data timestamp distribution in the *News* area. Our dataset spans from 1st May 2024 to 31th August 2024. By the time of evaluation, we inspect the knowledge

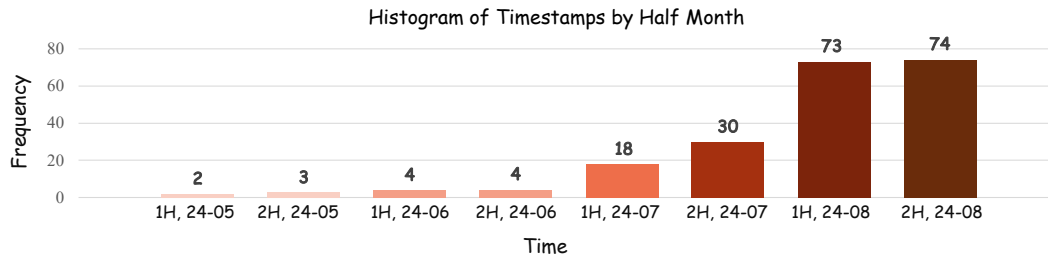


Figure 5: **Timestamps Distribution of Questions in the News Area.** All events of our collected data occurred after May 2024. The majority of data concentrates on August. This ensures the data captures only recent events, falling beyond the knowledge cutoff dates of LMMs. The False Premise subfield is not included, since it is infeasible to determine the timestamp for an event that never occurred.

cutoff dates of the closed-source models. Claude 3.5 Sonnet reports a knowledge cutoff of April 2024, while both GPT-4V and GPT-4o state they lack information from 2024. For open-source models, we examine their release dates and training data, confirming that none possess knowledge beyond May 2024. This temporal gap ensures the fairness of our evaluation, as the models’ performance solely reflects their multimodal search capabilities rather than pre-existing knowledge. We will update the *News* area if a new LMM’s training data may overlap with our collection period.

### 3 EXPERIMENT

In this section, we conduct a systematic evaluation of existing LMMs on MMSEARCH. We first introduce the experimental setup in Section 3.1. Then, we detail the quantitative results in Section 3.2 and narrate the error analysis in Section 3.3. Finally, we conduct an in-depth ablation study of scaling test-time compute versus scaling model size in Section E.3 in the supplementary material.

#### 3.1 EXPERIMENT SETUP

**Evaluation Models** We examine the performance of foundation models across three distinct categories on MMSEARCH: (a) *Commercial AI Search Engines*, represented by Perplexity (Perplexity). We test the pro version of Perplexity, which takes only the user query and image as input. Since SearchGPT (OpenAI, 2024c) has not been public yet, we do not test on it. (b) *Closed-source LMMs*, represented by models like GPT-4V (OpenAI, 2023c), GPT-4o (OpenAI, 2024b), and Claude 3.5 Sonnet (Anthropic, 2024), and (c) *Open-source LMMs*, featuring models such as LLaVA-OneVision-7B (Li et al., 2024b) (Qwen2-7B (Yang et al., 2024a)), LLaVA-OneVision-72B (Li et al., 2024b) (Qwen2-72B (Yang et al., 2024a)), LLaVA-NeXT-Interleave (Li et al., 2024c) (Qwen1.5-7B (Yang et al., 2024a)), InternVL2 (Chen et al., 2024d) (InternLM2.5-7B-Chat (Cai et al., 2024)), InternLM-XC2.5 (Zhang et al., 2024a) (InternLM2-7B (Cai et al., 2024)), Qwen2-VL-7B (Qwen Team, 2024) (Qwen2-7B (Yang et al., 2024a)), Qwen2-VL-72B (Qwen Team, 2024) (Qwen2-72B (Yang et al., 2024a)), mPlug-Owl3 (Ye et al., 2024) (Qwen2-7B (Yang et al., 2024a)), Idefics3 (Laurençon et al., 2024) (LLaMA3.1-7B-Instruct (AI@Meta, 2024)), and Mantis (Jiang et al., 2024b) (LLaMA3-7B (AI@Meta, 2024)). Note that the open-source LMMs’ sizes are 7B unless otherwise specified.

**Implementation Details.** We set the number of retrieved websites  $K$  as 8. We include two image input resolution settings. For the default settings, the longest edge of the input image is resized to match the largest resolution of the vision encoder of LMM. This ensures the image not to be cropped to multiple images and will only take up the minimum of tokens for image input. For any resolution settings, we input the image without resizing. More details are available in Section F.

#### 3.2 EXPERIMENTAL ANALYSIS

To thoroughly investigate the multimodal searching capabilities, we present the evaluation results of different models on MMSEARCH following the proposed step-wise evaluation strategy in Table 2

Table 2: **Evaluation Results of Final scores and Four Tasks in MMSEARCH.** We report the scores of news and knowledge areas and their average score in each task. Subscript *AnyRes* indicates original resolution image input; otherwise, low-resolution images were used. The highest scores for closed-source and open-source LMMs are marked in red and blue. For open-source LMMs, we adopt the models with 7B parameters unless otherwise specified.

Model	All			End-to-end			Requery			Rerank			Summarization		
	Avg	News	Know.	Avg	News	Know.	Avg	News	Know.	Avg	News	Know.	Avg	News	Know.
<i>Baselines</i>															
Human	69.2	69.6	68.1	68.2	68.6	67.1	43.7	45.0	40.1	85.7	87.3	81.2	72.8	71.4	76.7
<i>Commercial AI Search Engines</i>															
Perplexity Pro (Perplexity)	-	-	-	47.8	52.7	34.1	-	-	-	-	-	-	-	-	-
<i>Closed-source LMMs with MMSEARCH-ENGINE</i>															
Claude 3.5 Sonnet (Anthropic, 2024)	53.5	53.1	54.7	49.9	49.3	51.6	42.0	43.6	37.7	80.2	78.7	84.2	59.4	60.3	57.0
GPT-4V (OpenAI, 2023c)	55.0	55.0	55.3	52.1	52.2	51.9	45.7	49.2	35.8	79.3	76.9	86.1	57.4	56.7	59.4
GPT-4o (OpenAI, 2024b)	62.3	61.2	65.3	60.4	59.0	64.5	46.8	49.9	38.0	83.0	82.4	84.8	63.1	62.2	65.6
<i>Open-source LMMs with MMSEARCH-ENGINE</i>															
Mantis (Jiang et al., 2024b)	18.7	19.8	15.9	15.8	16.4	14.3	20.1	24.6	7.4	39.7	41.0	36.1	19.2	22.0	11.5
InternLM-XC2.5 (Zhang et al., 2024a)	22.2	22.8	20.5	22.9	23.6	20.8	25.0	24.3	27.0	0.0	0.0	0.0	37.7	38.6	35.0
InternLM-XC2.5 <sub>AnyRes</sub>	22.3	23.9	17.5	23.2	25.4	16.9	21.7	19.8	26.9	0.0	0.0	0.0	37.7	38.6	35.1
LLaVA-NeXT-Interleave (Li et al., 2024c)	28.3	29.2	25.6	23.0	23.8	20.5	26.2	30.7	13.5	55.3	58.6	46.2	42.5	40.0	49.3
mPlug-Owl3 (Ye et al., 2024)	32.1	34.8	24.4	24.6	28.1	14.9	32.6	36.7	21.2	74.3	73.5	76.6	45.6	45.6	45.4
mPlug-Owl3 <sub>AnyRes</sub>	33.9	35.5	29.3	27.3	29.4	21.2	31.8	36.1	19.9	74.5	72.9	79.1	43.9	43.6	44.6
InternVL2 (Chen et al., 2024d)	34.3	35.7	30.2	30.9	32.5	26.2	32.3	36.1	21.6	46.5	49.3	38.6	48.5	45.9	55.8
InternVL2 <sub>AnyRes</sub>	34.1	34.2	33.7	30.0	30.3	29.0	31.4	35.5	19.7	53.2	52.3	55.7	46.9	44.4	53.7
Idefics3 (Laurençon et al., 2024)	36.2	38.5	29.6	29.3	32.3	20.8	31.0	37.3	13.5	76.5	73.3	85.4	50.3	51.1	48.1
Idefics3 <sub>AnyRes</sub>	35.7	38.2	28.7	30.1	32.9	22.3	27.2	32.2	13.2	72.7	73.1	71.5	45.2	46.3	42.1
LLaVA-OneVision (Li et al., 2024b)	36.6	39.4	28.9	29.6	33.1	19.7	35.8	40.3	23.2	72.8	73.5	70.9	53.5	51.8	58.5
Qwen2-VL <sub>AnyRes</sub> (Qwen Team, 2024)	45.3	44.1	48.7	40.3	39.2	43.5	39.0	41.9	30.8	76.7	73.8	84.8	54.7	52.7	60.4
LLaVA-OneVision (72B)	50.1	50.1	50.2	44.9	45.1	44.1	42.9	45.9	34.3	82.2	80.5	86.7	61.4	59.1	67.7
Qwen2-VL <sub>AnyRes</sub> (72B)	52.7	52.0	54.5	49.1	48.8	49.8	44.7	47.2	37.6	76.7	72.9	87.3	59.6	57.5	65.7

and fourteen subfields in Table 3. We now provide a detailed discussion of notable findings and their implications for multimodal search capabilities.

**Any-resolution input only provides slight or no improvement.** Of the tested LMMs, four models, which are InternLM-XC2.5, InternVL2, mPlug-Owl3, and Idefic3, all support both low-resolution (LowRes) and any-resolution input (AnyRes). As one would expect, AnyRes input enables better OCR and perception of the image. However, we only observe slight or even no enhancement comparing the difference between the LowRes performance and its AnyRes counterpart. Take mPlug-Owl3 as an example, AnyRes input surpasses LowRes input on overall score by 1.8%, end-to-end score by 2.7%, and rerank on 0.2%. While it falls behind LowRes on requery by 0.8% and summarization by 1.7%. This suggests that the OCR and perception quality do not bottleneck the search performance. Rather, the suboptimal performance appears to stem from the LMMs’ inherent lack of robust search capabilities.

**Current LMMs still have significant shortcomings in requery and rerank.** Comparing the average score of the end-to-end task with that of the summarization task, we find that the summarization score consistently surpasses the end-to-end task by a large margin, both in the closed-source and open-source models. The minimum margin is 2.7% for GPT-4o, while the maximum is 23.9% for LLaVA-OneVision-7B. This discrepancy can be attributed to the differences in the tasks’ input quality. While the summarization task input always contains the answer, the end-to-end task’s third-round input quality depends on the model’s requery and rerank quality in previous rounds. The magnitude of this performance gap reflects the disparity between a model’s summarization ability and its capacity for requery and rerank tasks. The larger the difference, the larger the capability gap. Observing the result, we find that this gap of most open-sourced models exceeds 14%, while the closed-sourced models are all below 10%. This suggests all current LMMs needs improvement of their requery and rerank ability, especially for open-source models. Mantis is one exception of open-source models with a margin of only 3.4%. This means its poor summarization capability bottlenecks its end-to-end performance. Qwen2-VL-72B’s 10.5% gap, also falling below 14%, highlights its superiority among other open-source LMMs.



Table 3: **Evaluation Results on Different Subfields in MMSEARCH.** SPO: Traditional Sports; ESP: E-Sports; ENT: Entertainment; GEN: General News; PAP: Paper; TEC: Technology; FIN: Finance; FAL: False Premise; ART: Arts; ARC: Architecture; AST: Astronomy; ANI: Anime; AUT: Auto; FAS: Fashion. The highest scores for closed-source and open-source LMMs are marked in red and blue. For open-source LMMs, we adopt the models with 7B sizes unless otherwise specified.

Model	All		News									Knowledge					
	Avg		SPO	ESP	ENT	GEN	PAP	TEC	FIN	FAL	Avg	ART	ARC	AST	ANI	AUT	FAS
<i>Closed-source LMMs with MMSEARCH-ENGINE</i>																	
Claude 3.5 Sonnet (Anthropic, 2024)	53.5	53.0	37.4	50.2	63.6	49.2	52.8	67.7	43.4	63.1	54.7	59.0	42.9	70.9	56.5	60.6	36.4
GPT-4V (OpenAI, 2023c)	55.0	54.9	49.3	48.5	67.4	43.5	37.3	64.2	59.1	90.2	55.3	54.6	45.6	65.1	63.0	52.2	56.2
GPT-4o (OpenAI, 2024b)	62.3	61.2	63.7	61.2	72.3	51.3	48.6	68.6	60.0	76.8	65.3	73.8	52.0	57.6	76.8	68.4	55.8
<i>Open-source LMMs with MMSEARCH-ENGINE</i>																	
InternLM-XC2.5 (Zhang et al., 2024a)	22.2	22.8	22.6	13.6	28.9	23.5	15.5	27.8	32.9	3.1	20.5	29.2	24.0	20.8	2.0	20.4	11.8
InternLM-XC2.5 <sub>AnyRes</sub>	22.2	23.9	25.2	12.8	30.6	21.0	18.7	31.4	32.8	13.9	17.6	22.9	25.6	10.8	2.0	28.1	4.7
Mantis (Jiang et al., 2024b)	18.8	19.8	17.4	9.7	25.3	19.6	22.1	35.3	18.2	6.2	15.9	23.0	17.5	12.5	14.8	12.7	3.3
LLaVA-NeXT-Interleave (Li et al., 2024c)	28.3	29.3	23.3	18.5	41.9	33.6	26.2	27.3	28.0	14.7	25.6	31.1	22.7	49.7	16.5	19.7	7.1
InternVL2 (Chen et al., 2024d)	34.3	35.7	38.7	20.7	46.8	33.2	24.6	44.1	44.5	27.2	30.1	41.7	33.5	26.9	15.4	36.8	7.9
InternVL2 <sub>AnyRes</sub>	34.0	34.2	32.0	26.6	45.4	35.0	18.7	23.7	43.1	36.3	33.7	49.5	26.4	36.5	18.3	35.2	15.5
mPlug-Owl3 (Ye et al., 2024)	32.1	34.8	30.0	20.6	46.7	33.7	22.3	30.4	45.7	41.5	24.4	28.7	24.6	33.5	18.3	24.2	10.5
mPlug-Owl3 <sub>AnyRes</sub>	33.9	35.5	42.0	27.1	49.7	34.4	24.5	30.4	41.3	18.9	29.3	34.4	20.3	40.9	18.8	32.8	24.4
Idefics3 (Laurençon et al., 2024)	36.2	38.5	43.5	24.2	50.1	41.7	27.5	36.2	38.8	37.4	29.6	34.2	32.1	32.1	29.6	32.7	9.1
Idefics3 <sub>AnyRes</sub>	35.7	38.2	40.5	27.3	48.0	41.4	33.1	36.9	42.1	17.4	28.8	40.5	26.2	27.1	14.2	40.8	7.3
LLaVA-OneVision (Li et al., 2024b)	36.6	39.4	31.8	27.1	50.0	38.5	31.4	52.1	47.4	23.5	28.9	38.0	26.4	34.3	21.0	29.6	11.2
Qwen2-VL <sub>AnyRes</sub> (Qwen Team, 2024)	45.3	44.1	47.8	33.9	49.4	45.8	45.6	38.1	49.8	31.1	48.7	70.0	32.7	43.5	42.5	54.4	23.3
LLaVA-OneVision (72B)	50.1	52.3	53.2	45.4	62.1	45.6	41.5	64.8	47.8	37.2	44.5	63.4	42.8	52.4	24.3	70.0	31.7
Qwen2-VL <sub>AnyRes</sub> (72B)	52.7	52.9	50.0	46.3	66.2	36.3	55.0	57.2	56.6	52.0	52.1	58.0	45.8	65.4	45.7	68.6	41.9

**Closed-source LMMs are better-performed than open-sourced LMMs on overall performance.** For the final score, closed-source LMMs consistently outperform the open-source LMMs. GPT-4o achieves the highest overall score of 62.3%, demonstrating superior zero-shot multimodal search capabilities. While Qwen2-VL-72B leads among open-source models, it still lags behind GPT-4o by 9.6%. The performance gap widens to 11.3% on the most challenging end-to-end task and further expands to 20.1% for 7B open-source LMMs. These significant disparities highlight substantial room for improvement in open-source models.

**SoTA LMMs with our MMSEARCH-ENGINE surpass commercial AI search engines in the end-to-end task.** We also evaluate the pro version of Perplexity (Perplexity), a prominent commercial AI search engine that accepts both image and text queries, on our dataset. Perplexity pro can accept both image and text in the user query. Surprisingly, although Perplexity also leverages SoTA LMMs like GPT-4o and Claude 3.5 Sonnet, it largely underperforms MMSEARCH-ENGINE equipped with the same model in the end-to-end task. Even more remarkably, MMSEARCH-ENGINE can even surpass Perplexity with Qwen2-VL-72B, an open-source LMM. *This suggests that our MMSEARCH-ENGINE provides a better open-source plan for multimodal AI search engine.* The performance gap validates MMSEARCH-ENGINE’s design effectiveness and highlights the value of testing various LMMs within our pipeline, since the pipeline can indeed achieve remarkable performance when using powerful LMMs. Upon investigating Perplexity’s sub-optimal performance, we discovered that it appears to utilize only a rudimentary image search algorithm, if any. This limitation leads to poor identification of the key objects in the image and failure to retrieve relevant information. Our findings underscore the effectiveness of MMSEARCH-ENGINE’s design, particularly the incorporation of a robust image search step, which plays a crucial role in accurately recognizing important information from the input image.

### 3.3 ERROR ANALYSIS

To investigate the limitations of current LMM search capabilities, we conducted a comprehensive analysis of error types observed in our evaluation. Our proposed step-wise evaluation strategy enables analysis of failure modes for each core search step, complementing the end-to-end assessment. This analysis encompasses the entire benchmark. We first examine the end-to-end error types for both the best-performing closed-source model (GPT-4o) and open-source model (Qwen2-VL-7B).

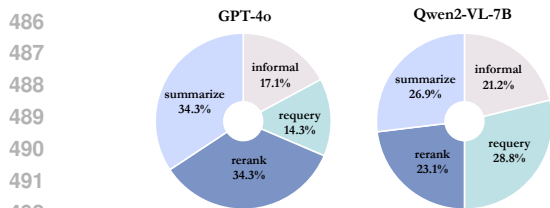


Figure 6: Distribution of Error types of GPT-4o (OpenAI, 2024b) and Qwen2-VL-7B (Qwen Team, 2024) in the end-to-end task.

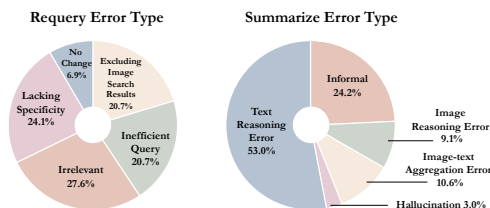


Figure 7: Distribution of Error types of Qwen2-VL-7B (Qwen Team, 2024) in the requery and summarization task respectively.

To better understand the failure cases, we then identify distinct error types in the requery and summarization task, which requires open-ended generation. We quantify these error types for a systematic understanding of current LMM limitations and point out critical areas for improvement.

### 3.3.1 ERROR ANALYSIS OF END-TO-END TASK

In this section, we are trying to answer the question: *Which step does LMM make a mistake in the end-to-end evaluation?* In Fig. 6, we showcase the statistics of different error types occurring in GPT-4o and Qwen2-VL-7B. We define the following four error categories: (i) *requery*, where the model requery is incorrect, and leads to all retrieved websites being invalid; (ii) *rerank*, where the model selects a website without a correct answer; (iii) *summarization*, where the full website content contains the information of correct answer, but the model fails to extract it; (iv) *informal*, the output format deviates from the prompt specifications. As shown in the figure, GPT-4o’s primary error sources are rerank and summarization errors, while requery and informal errors account for approximately half the frequency of the main error causes. This suggests that GPT-4o’s limitations lie primarily in information source ranking and multimodal information integration. As for Qwen2-VL, all four error types occur with similar frequency. The rise of the informal error portion may be attributed to the model’s inferior instruction-following ability. Besides, it should be noted that the requery task demands advanced comprehension and key image information extraction ability. This task seldom appears in the training data of current LMMs. The prevalence of this error type in Qwen2-VL may indicate that it fails to generalize to adequately address this complex task.

### 3.3.2 ERROR ANALYSIS OF REUQUERY AND SUMMARIZATION TASK

To better understand how open-source LMM makes the mistake, we dive into the requery and summarization task to find out the error patterns of Qwen2-VL-7B. We particularly select the two tasks requiring open-ended generation, which provides more information to identify the error.

We define five types of error for both the requery and summarization tasks in Section E.4 and showcase the statistics in Fig. 7. The requery errors suggest that LMM often fails to fully understand the requery task and aggregate all available information. Besides, inefficient query error indicates that LMM has no clue about the real working scenario and query principles of search engines. On the other hand, the five types of summarization errors reflect that current LMMs still cannot correctly extract the given multimodal information to answer the query. The ability of content understanding still requires further enhancement.

## 4 CONCLUSION

In this paper, we investigate the potential of LMMs as multimodal AI search engines. We first design MMSEARCH-ENGINE, a streamlined pipeline, enabling zero-shot LMMs to perform multimodal searches. To comprehensively assess the search capabilities, we introduce MMSEARCH, a benchmark comprising 300 queries across 14 subfields. Our evaluation methodology analyzes LMM search abilities step-by-step, facilitating a deeper understanding of their limitations. Using MMSEARCH-ENGINE, we evaluate various closed-source and open-source LMMs, revealing that current models still fall short of human-level search proficiency. Through thorough error analysis, we identify specific patterns of failure in key search process steps, providing valuable insights for future improvements in LMM search ability.

540 REPRODUCIBILITY STATEMENT

541  
542 We provide the demo code of MMSEARCH-ENGINE, which can be also used for inference, in  
543 the *code* directory in the supplementary material. We also include the end-to-end data in the *data*  
544 directory in the supplementary material, along with the loading script. Implementation details of  
545 inference settings are demonstrated in Section 3.1 and Section F.

547 REFERENCES

- 548  
549 Google lens. URL <https://lens.google.com>. Web interface available at  
550 <https://images.google.com>.
- 551 AI@Meta. Llama 3 model card, 2024. URL [https://github.com/meta-llama/llama3/  
552 blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 553  
554 Anthropic. Claude-3.5. <https://www.anthropic.com/news/claude-3-5-sonnet>,  
555 2024.
- 556 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to  
557 retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- 558  
559 Oriol Barbany, Michael Huang, Xinliang Zhu, and Arnab Dhua. Leveraging large language models  
560 for multimodal search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
561 Pattern Recognition*, pp. 1201–1210, 2024.
- 562 Patrice B chard and Orlando Marquez Ayala. Reducing hallucination in structured outputs via  
563 retrieval-augmented generation. *arXiv preprint arXiv:2404.08189*, 2024.
- 564  
565 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Milli-  
566 can, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al.  
567 Improving language models by retrieving from trillions of tokens. In *International conference on  
568 machine learning*, pp. 2206–2240. PMLR, 2022.
- 569 Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine.  
570 *Computer networks and ISDN systems*, 30(1-7):107–117, 1998.
- 571  
572 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui  
573 Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye  
574 Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting  
575 Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaying Li, Jingwen Li, Linyang Li,  
576 Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun  
577 Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang  
578 Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song,  
579 Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang,  
580 Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong,  
581 Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia  
582 Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo  
583 Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, Hui  
584 Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng Zou,  
Xipeng Qiu, Yu Qiao, and Dahua Lin. Internlm2 technical report, 2024.
- 585  
586 Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. Rq-rag:  
587 Learning to refine queries for retrieval augmented generation. *arXiv preprint arXiv:2404.00610*,  
2024.
- 588  
589 Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang,  
590 Yu Qiao, Tong Lu, et al. Videollm: Modeling video sequence with large language models. *arXiv  
591 preprint arXiv:2305.13292*, 2023a.
- 592  
593 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding:  
Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge dis-  
tillation, 2024a.

- 594 Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in  
595 retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
596 volume 38, pp. 17754–17762, 2024b.
- 597 Jun Chen, Deyao Zhu<sup>1</sup> Xiaoqian Shen<sup>1</sup> Xiang Li, Zechun Liu<sup>2</sup> Pengchuan Zhang, Raghuraman  
598 Krishnamoorthi<sup>2</sup> Vikas Chandra<sup>2</sup> Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: Large  
599 language model as a unified interface for vision-language multi-task learning. *arXiv preprint*  
600 *arXiv:2310.09478*, 2023b.
- 601 Xingyu Chen, Zihan Zhao, Lu Chen, Danyang Zhang, Jiabao Ji, Ao Luo, Yuxuan Xiong, and  
602 Kai Yu. Websrc: A dataset for web-based structural reading comprehension. *arXiv preprint*  
603 *arXiv:2101.09465*, 2021.
- 604 Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning Liu, Wenwei Zhang, Kai Chen, and Feng Zhao.  
605 Mindsearch: Mimicking human minds elicits deep ai searcher. *arXiv preprint arXiv:2407.20183*,  
606 2024c.
- 607 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,  
608 Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to com-  
609 mercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024d.
- 610 Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei,  
611 Songyang Zhang, Haodong Duan, Maosong Cao, et al. Internlm-xcomposer2: Mastering free-  
612 form text-image composition and comprehension in vision-language large model. *arXiv preprint*  
613 *arXiv:2401.16420*, 2024.
- 614 Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and  
615 Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In  
616 *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,  
617 pp. 6491–6501, 2024.
- 618 Brian J Fogg, Jonathan Marshall, Othman Laraki, Alex Osipovich, Chris Varma, Nicholas Fang,  
619 Jyoti Paul, Akshay Rangnekar, John Shon, Preeti Swani, et al. What makes web sites credible? a  
620 report on a large quantitative study. In *Proceedings of the SIGCHI conference on Human factors*  
621 *in computing systems*, pp. 61–68, 2001.
- 622 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu  
623 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation  
624 benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- 625 Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu  
626 Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evalua-  
627 tion benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
- 628 Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong,  
629 Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal  
630 large language model. *arXiv preprint arXiv:2312.11370*, 2023a.
- 631 Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu,  
632 Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient  
633 visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023b.
- 634 Peng Gao, Renrui Zhang, Chris Liu, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie  
635 Geng, Ziyi Lin, Peng Jin, et al. Sphinx-x: Scaling data and parameters for a family of multi-modal  
636 large language models. *ICML 2024*, 2024.
- 637 Google Gemini Team. Gemini: a family of highly capable multimodal models. *arXiv preprint*  
638 *arXiv:2312.11805*, 2023.
- 639 Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen,  
640 Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud  
641 with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint*  
642 *arXiv:2309.00615*, 2023.



- 648 Ziyu Guo, Renrui Zhang, Hao Chen, Jialin Gao, Peng Gao, Hongsheng Li, and Pheng-Ann  
649 Heng. Sciverse. [https://sciverse-cuhk.  
650 github.io/](https://sciverse-cuhk.github.io).  
651
- 652 Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Chengzhuo Tong, Peng Gao, Chunyuan Li, and Pheng-  
653 Ann Heng. Sam2point: Segment any 3d as videos in zero-shot and promptable manners. *arXiv  
654 preprint arXiv:2408.16768*, 2024b.
- 655 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented  
656 language model pre-training. In *International conference on machine learning*, pp. 3929–3938.  
657 PMLR, 2020.  
658
- 659 Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu,  
660 Song Wen, Ziyu Guo, et al. Imagebind-llm: Multi-modality instruction tuning. *arXiv preprint  
661 arXiv:2309.03905*, 2023.
- 662 Qiuxiang He, Guoping Huang, Qu Cui, Li Li, and Lemao Liu. Fast and accurate neural machine  
663 translation with translation memory. In *Proceedings of the 59th Annual Meeting of the Association  
664 for Computational Linguistics and the 11th International Joint Conference on Natural Language  
665 Processing (Volume 1: Long Papers)*, pp. 3170–3180, 2021.  
666
- 667 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong  
668 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language  
669 models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*,  
670 2023.
- 671 Bowen Jiang, Yangxinyu Xie, Xiaomeng Wang, Weijie J Su, Camillo J Taylor, and Tanwi  
672 Mallick. Multi-modal and multi-agent systems meet rationality: A survey. *arXiv preprint  
673 arXiv:2406.00252*, 2024a.  
674
- 675 Dongfu Jiang, Xuan He, Huaye Zeng, Con Wei, Max Ku, Qian Liu, and Wenhua Chen. Mantis:  
676 Interleaved multi-image instruction tuning. *arXiv preprint arXiv:2405.01483*, 2024b.
- 677 Dongzhi Jiang, Guanglu Song, Xiaoshi Wu, Renrui Zhang, Dazhong Shen, Zhuofan Zong, Yu Liu,  
678 and Hongsheng Li. Comat: Aligning text-to-image diffusion model with image-to-text concept  
679 matching. *arXiv preprint arXiv:2404.03653*, 2024c.  
680
- 681 Nick Kanopoulos, Nagesh Vasanthavada, and Robert L Baker. Design of an image edge detection  
682 filter using the sobel operator. *IEEE Journal of solid-state circuits*, 23(2):358–367, 1988.
- 683 Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better under-  
684 standing vision-language models: insights and future directions., 2024.  
685
- 686 Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu,  
687 and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild.  
688 <https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms/>, 2024a.
- 689 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei  
690 Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint  
691 arXiv:2408.03326*, 2024b.  
692
- 693 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.  
694 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv  
695 preprint arXiv:2407.07895*, 2024c.
- 696 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-  
697 training for unified vision-language understanding and generation. In *International Conference  
698 on Machine Learning*, pp. 12888–12900. PMLR, 2022.  
699
- 700 KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang,  
701 and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*,  
2023.

- 702 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to  
703 solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024d.
- 704
- 705 Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi  
706 Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for  
707 multi-modal large language models. *ECCV 2024*, 2023.
- 708 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,  
709 2023a.
- 710 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.  
711 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL [https://  
712 llava-vl.github.io/blog/2024-01-30-llava-next/](https://llava-vl.github.io/blog/2024-01-30-llava-next/).
- 713
- 714 Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang  
715 Yue. Visualwebbench: How far have multimodal llms evolved in web page understanding and  
716 grounding? *arXiv preprint arXiv:2404.05955*, 2024b.
- 717 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan,  
718 Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around  
719 player? *arXiv preprint arXiv:2307.06281*, 2023b.
- 720
- 721 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,  
722 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for  
723 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521,  
724 2022.
- 725 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun yue Li, Hannaneh Hajishirzi, Hao Cheng, Kai-  
726 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating math reasoning in visual  
727 contexts with gpt-4v, bard, and other large multimodal models. *ArXiv*, abs/2310.02255, 2023.
- 728
- 729 Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-  
730 turn dialogue. *arXiv preprint arXiv:2402.05930*, 2024.
- 731 Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. Query rewriting for retrieval-  
732 augmented large language models. *arXiv preprint arXiv:2305.14283*, 2023.
- 733
- 734 OpenAI. Chatgpt. <https://chat.openai.com>, 2023a.
- 735
- 736 OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023b.
- 737
- 738 OpenAI. GPT-4V(ision) system card, 2023c. URL [https://openai.com/research/  
gpt-4v-system-card](https://openai.com/research/gpt-4v-system-card).
- 739
- 740 OpenAI. Openai o1. [Online], 2024a. [https://openai.com/index/  
learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/).
- 741
- 742 OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o/>, 2024b.
- 743
- 744 OpenAI. Searchgpt. [Online], 2024c. [https://openai.com/index/  
searchgpt-prototype/](https://openai.com/index/searchgpt-prototype/).
- 745
- 746 Perplexity. Perplexity.ai. [Online]. <https://www.perplexity.ai/>.
- 747
- 748 Qwen Team. Qwen2-vl. 2024.
- 749
- 750 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
751 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
752 Sutskever. Learning transferable visual models from natural language supervision. In *Internat-  
753 ional Conference on Machine Learning*, 2021. URL [https://api.semanticscholar.  
org/CorpusID:231591445](https://api.semanticscholar.org/CorpusID:231591445).
- 754
- 755 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.

- 756 Elizabeth Silience, Pam Briggs, Lesley Fishwick, and Peter Harris. Trust and mistrust of online  
757 health sites. In *Proceedings of the SIGCHI conference on Human factors in computing systems*,  
758 pp. 663–670, 2004.
- 759 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide  
760 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF*  
761 *Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.
- 762 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
763 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
764 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 765 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
766 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
767 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 768 Fei Wang, Xingyu Fu, James Y Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma,  
769 Nan Xu, Wenxuan Zhou, Kai Zhang, et al. Muirbench: A comprehensive benchmark for robust  
770 multi-image understanding. *arXiv preprint arXiv:2406.09411*, 2024.
- 771 Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. Large multimodal agents: A  
772 survey. *arXiv preprint arXiv:2402.15116*, 2024.
- 773 Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm:  
774 Empowering large language models to understand point clouds. *arXiv preprint arXiv:2308.16911*,  
775 2023.
- 776 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
777 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
778 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
779 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng  
780 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
781 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
782 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
783 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
784 technical report. *arXiv preprint arXiv:2407.10671*, 2024a.
- 785 Senqiao Yang, Jiaming Liu, Ray Zhang, Mingjie Pan, Zoey Guo, Xiaoqi Li, Zehui Chen, Peng  
786 Gao, Yandong Guo, and Shanghang Zhang. Lidar-llm: Exploring the potential of large language  
787 models for 3d lidar understanding. *arXiv preprint arXiv:2312.14074*, 2023.
- 788 Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary,  
789 Rongze Daniel Gui, Ziran Will Jiang, Ziyu Jiang, et al. Crag-comprehensive rag benchmark.  
790 *arXiv preprint arXiv:2406.04744*, 2024b.
- 791 Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and  
792 Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large  
793 language models. *arXiv preprint arXiv:2408.04840*, 2024.
- 794 Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen  
795 Hu, Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng  
796 Tian, Qi Qian, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language  
797 models with multimodality, 2023a.
- 798 Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei  
799 Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with  
800 modality collaboration, 2023b.
- 801 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,  
802 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *ArXiv*,  
803 abs/2308.02490, 2023.
- 804 Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language  
805 model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.

810 Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong  
811 Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language  
812 model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024a.  
813

814 Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng  
815 Gao, and Yu Qiao. LLaMA-adapter: Efficient fine-tuning of large language models with zero-  
816 initialized attention. In *The Twelfth International Conference on Learning Representations*,  
817 2024b. URL <https://openreview.net/forum?id=d4UiXAHN2W>.

818 Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou,  
819 Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the  
820 diagrams in visual math problems? *ECCV 2024*, 2024c.

821 Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Yichi Zhang, Ziyu Guo, Chengzhuo Tong, Jiaming Liu,  
822 Aojun Zhou, Bin Wei, Shanghang Zhang, et al. Mavis: Mathematical visual instruction tuning.  
823 *arXiv preprint arXiv:2407.08739*, 2024d.  
824

825 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: En-  
826 hancing vision-language understanding with advanced large language models. *arXiv preprint*  
827 *arXiv:2304.10592*, 2023.

828 Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu Song, Hao Shao, Dongzhi Jiang, Hongsheng  
829 Li, and Yu Liu. Mova: Adapting mixture of vision experts to multimodal context. *arXiv preprint*  
830 *arXiv:2404.13046*, 2024.  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



## SUPPLEMENTARY MATERIAL OVERVIEW

- Section **A**: Related work.
- Section **B**: Automated data curation pipeline.
- Section **C**: Additional dataset details.
- Section **D**: Future direction.
- Section **E**: Additional experiments and analysis.
- Section **F**: Additional experimental details.
- Section **G**: More dataset details.
- Section **H**: Qualitative examples.

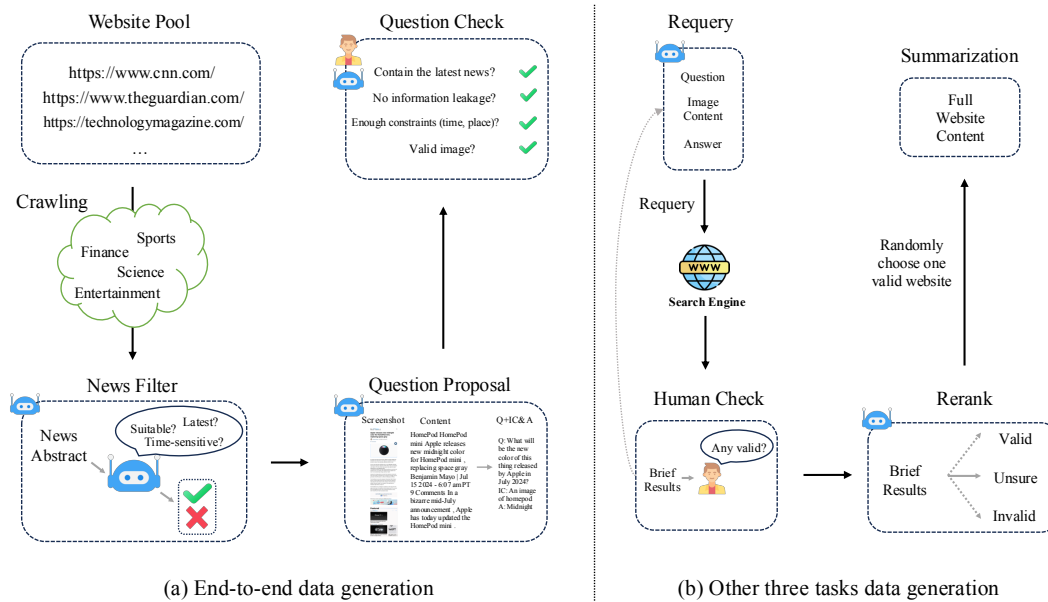
## A RELATED WORK

**Large Multimodal Models.** Recently, multimodal models (Radford et al., 2021; Li et al., 2022; OpenAI, 2023c; Rombach et al., 2022; Jiang et al., 2024c) has gained unparalleled attention. Building on the success of Large Language Models (LLMs) (Touvron et al., 2023a;b) and large-scale vision models (Radford et al., 2021), Large Multimodal Models (LMMs) are gaining prominence across diverse domains. These models extend LLMs to handle tasks involving various modalities, including mainstream 2D image processing (Liu et al., 2023a; Zhu et al., 2023; Lin et al., 2023; Gao et al., 2023b), as well as 3D point clouds (Xu et al., 2023; Guo et al., 2023; 2024b), and videos (Li et al., 2023; Chen et al., 2023a; Zhang et al., 2023; Fu et al., 2024). Among these LMMs, OpenAI’s GPT-4o (OpenAI, 2024b) and Anthropic’s Claude 3.5 Sonnet (Anthropic, 2024) demonstrate outstanding visual reasoning and comprehension capability, setting new standards in multi-modal performance. However, their closed-source nature limits broader adoption and development. In contrast, another research trajectory focuses on open-source LMMs for the community. Pioneering works like LLaVA (Liu et al., 2023a; 2024a; Li et al., 2024c;b), LLaMA-Adapter (Zhang et al., 2024b; Gao et al., 2023b), and MiniGPT-4 (Zhu et al., 2023; Chen et al., 2023b) incorporate a frozen CLIP (Radford et al., 2021) model for image encoding and integrate visual information into LLM for multi-modal instruction tuning. Later, works such as mPLUG-Owl (Ye et al., 2023a;b; 2024), SPHINX (Gao et al., 2024; Lin et al., 2023), and InternLM-XComposer (Dong et al., 2024) further advanced the field by incorporating diverse visual instruction tuning data and generalizing to more scenarios. More recent developments in the field have taken diverse directions. For example, several studies (Zong et al., 2024; Tong et al., 2024) explore multiple vision encoders design. Meanwhile, other works (Liu et al., 2024a; Chen et al., 2024d; Qwen Team, 2024) incorporate high-resolution image input. Multi-image instruction data (Li et al., 2024c; Jiang et al., 2024b) is also integrated to enable perception across multiple images. While various benchmarks, both in the general (Fu et al., 2023; Liu et al., 2023b; Yu et al., 2023) and expert (Zhang et al., 2024c; Lu et al., 2023; 2022) domain, has been proposed, the potential of LMM to function as a multimodal search engine remains largely unexplored. To this end, we introduce the MMSEARCH benchmark, which evaluates LMMs’ zero-shot abilities of multimodal search, offering valuable insights for future research.

**Large models with Retrieval Augmented Generation (RAG).** RAG (Retrieval-Augmented Generation) is an effective strategy for enhancing model knowledge by retrieving relevant information from external sources (Fan et al., 2024). RAG has been leveraged in various scenarios including knowledge-intensive question answering (Borgeaud et al., 2022; Guu et al., 2020), machine translation (He et al., 2021), and hallucination elimination (Bécharde & Ayala, 2024). Current works has focused on improving specific aspects of RAG. RG-RAG (Chan et al., 2024) proposes to refine the query for retrieval by decomposition and disambiguation. Self-RAG (Asai et al., 2023) incorporates the self-reflection of LLM to enhance the generation quality. The AI search engine could be viewed as a form of RAG with the Internet serving as the external knowledge source. Recently, MindSearch (Chen et al., 2024c) proposes an AI search engine framework to simulate the human minds in web information seeking. Meanwhile, multiple benchmarks of RAG (Yang et al., 2024b; Chen et al., 2024b) have been introduced to comprehensively evaluate a RAG system. However, both the current AI search engine and RAG benchmark are limited to the text-only setting, leaving the multimodal search engine and evaluation largely unexplored. To bridge this gap, we introduce

MMSEARCH-ENGINE and MMSEARCH, a multimodal AI search engine pipeline and dataset designed to evaluate various multimodal scenarios.

## B AUTOMATED DATA CURATION PIPELINE



**Figure 8: Illustration of the Data Curation Pipeline.** We first collect the end-to-end data by crawling inside the website pool and prompting LMMs to raise questions based on the content. Then we further prompt the LMM to generate annotation for other three tasks. The human check is optional for the end-to-end data generation but compulsory for other three tasks data generation.

Now we introduce our automated/semi-automated data curation pipeline. The figure is shown in Fig. 8. We first define a website pool and a model pool. The website pool contains general news websites like CNN and expertise websites like Arxiv. The model pool contains state-of-art models for the data curation pipeline to guarantee diversity and fairness.

- 1. End-to-end data curation.** We employ a web crawler to obtain all the sub-websites published later than a specific date. However, not all the websites are suitable for raising a question to test the LMMs’ searching capability. For example, some websites do not contain any recent news, while some websites’ contents are difficult to convert into a question with a definite answer. Therefore, we randomly choose a model from the model pool to serve as a news filter, prompting it to filter valid websites by providing some few-shot examples. Next, we provide the text contents and screenshots of the valid websites to a model from the model pool. The model is asked to raise several questions based on the website content. It is encouraged to raise questions unable to be answered only by text. As for the question with an image, the model is asked to briefly describe the image content, and we will later use the description to search in Bing and obtain the first result image. Finally, we apply the quality check of the generated questions and their corresponding images either by human or a model from the model pool.
- 2. Rerank data generation.** We provide the generated requery to the search engine and retrieve  $K$  websites for rerank. Again, we provide the question, the question image content, and the answer to a model and ask to categorize each website into valid, unsure or invalid.
- 3. Summarization data generation.** We randomly choose one website from the websites marked as valid in the last step and obtain its full content.

**Table 4: Evaluation Results of Different Complexity in MMSEARCH.** We categorize all the questions to three levels of complexity.

Model	All			End-to-end			Requery			Rerank			Summarization		
	Easy	Middle	Hard	Easy	Middle	Hard	Easy	Middle	Hard	Easy	Middle	Hard	Easy	Middle	Hard
<i>Closed-source LMMs with MMSEARCH-ENGINE</i>															
Claude 3.5 Sonnet (Anthropic, 2024)	56.2	53.6	48.9	53.0	49.8	45.0	49.3	37.6	35.0	80.3	84.5	75.3	59.8	59.9	58.3
GPT-4V (OpenAI, 2023c)	55.3	56.4	53.1	52.7	52.9	50.2	52.7	41.4	38.8	76.5	82.8	80.2	55.0	63.9	54.3
GPT-4o (OpenAI, 2024b)	63.1	63.9	59.2	61.0	62.3	57.4	54.3	41.2	40.6	82.2	85.6	81.5	64.0	65.5	59.0
<i>Open-source LMMs with MMSEARCH-ENGINE</i>															
Mantis (Jiang et al., 2024b)	22.8	18.0	13.0	19.4	15.9	10.0	28.5	12.6	14.5	40.5	43.1	34.6	28.2	11.5	12.9
InternLM-XC2.5 (Zhang et al., 2024a)	26.6	19.3	18.1	28.5	20.0	16.7	24.3	26.0	25.0	0.0	0.0	0.0	39.5	30.3	42.7
InternLM-XC2.5 <sub>AnyRes</sub>	25.6	18.1	21.2	27.5	18.3	21.4	20.2	24.6	20.9	0.0	0.0	0.0	39.5	31.5	41.4
LLaVA-NeXT-Interleave (Li et al., 2024c)	33.2	25.6	23.3	28.5	20.2	16.9	31.0	22.1	22.7	61.4	51.7	49.4	41.5	41.6	45.0
mPlug-Owl3 (Ye et al., 2024)	37.1	28.5	27.7	30.5	20.6	19.5	40.8	28.3	23.9	77.7	73.0	70.4	44.8	44.1	48.4
mPlug-Owl3 <sub>AnyRes</sub>	37.5	32.8	29.0	31.4	26.4	21.5	38.6	26.9	26.0	76.5	77.6	67.9	44.0	39.6	48.3
InternVL2 (Chen et al., 2024d)	38.7	30.9	30.5	35.0	28.3	26.8	39.5	25.3	27.9	57.6	38.5	37.0	47.4	46.2	52.8
InternVL2 <sub>AnyRes</sub>	38.0	30.0	32.0	34.6	24.5	28.2	38.8	24.5	26.5	53.8	59.8	45.1	46.7	43.9	50.2
Idefics3 (Laureçon et al., 2024)	43.1	30.5	31.0	38.0	21.3	23.5	40.9	23.3	23.3	76.5	83.3	69.1	48.9	50.0	53.1
Idefics3 <sub>AnyRes</sub>	43.2	31.2	28.4	38.4	24.1	23.0	37.6	20.3	17.6	76.5	74.7	64.2	48.4	46.6	38.4
LLaVA-OneVision (Li et al., 2024b)	42.9	32.3	31.1	36.8	24.6	23.3	46.5	28.9	25.8	78.0	67.8	69.8	51.5	56.7	53.6
Qwen2-VL <sub>AnyRes</sub> (Qwen Team, 2024)	45.6	44.9	45.3	41.2	38.9	40.4	46.5	34.0	32.2	73.5	83.3	74.7	50.3	56.8	59.6
LLaVA-OneVision (72B)	52.2	48.0	49.0	47.6	41.6	44.0	51.0	39.6	33.2	79.5	85.1	83.3	60.3	63.7	60.6
Qwen2-VL <sub>AnyRes</sub> (72B)	53.7	53.4	50.2	50.7	48.6	46.9	52.3	40.0	37.1	72.3	82.2	77.8	58.5	67.0	53.6

Notably, human check is a must for the requery data generation process. There should be at least one valid website to guarantee the effectiveness of the generated requery. Only after human check of this step, the quality of rerank and summarization data generation is assured.

## C ADDITIONAL DATASET DETAILS

We manually annotate the complexity of the data based on the difficulty of the three steps in MMSearch-Engine:

- 1. Requery difficulty.** This concerns the complexity of transforming the original question into an effective search query. Complex cases arise when the question references image content, requiring the LMM to first analyze the visual information and then synthesize it with the text question into a coherent search query. For instance, when a user asks about a landmark shown in an image, the LMM must first identify the landmark through image search and then incorporate this information into a text query about the landmark’s specific attributes.
- 2. Rerank difficulty.** This dimension evaluates the challenge of identifying and prioritizing relevant search results. The difficulty primarily scales with the information scarcity. If there are only very limited websites containing useful information, it is more difficult to successfully retrieve and choose the website.
- 3. Summarization difficulty.** This aspect involves both information synthesis and multi-modal reasoning challenges. In cases requiring synthesis, answers cannot be derived from a single source sentence - the LMM must integrate information scattered across different parts of the website. For example, comparing event frequencies across locations (like concert counts between cities) requires gathering and analyzing distributed data. Additionally, some questions demand analysis of both textual and visual website content, sometimes necessitating comparison with input images with images in the website.

Based on these criteria, we have categorized all questions into three difficulty levels, with the following distribution: hard (28%), medium (27.7%), and easy (44.3%). We provide the evaluation results grouped by the difficulty levels in Table 4.

## D FUTURE DIRECTION

Our proposed MMSEARCH-ENGINE can be enhanced through interactive user feedback loops. When the model produces an incorrect answer, users can identify the specific step where the error occurred and prompt the model to reconsider its reasoning. This iterative process allows for guided model refinement until accurate results are achieved. We conducted preliminary experiments on GPT-4o with this approach on three test cases. The results demonstrated that the LMM successfully interpreted user feedback and appropriately adjusted its responses based on the provided guidance. The model’s ability to understand the user input and revise its reasoning suggests significant potential for improving task accuracy.

## E ADDITIONAL EXPERIMENTS AND ANALYSIS

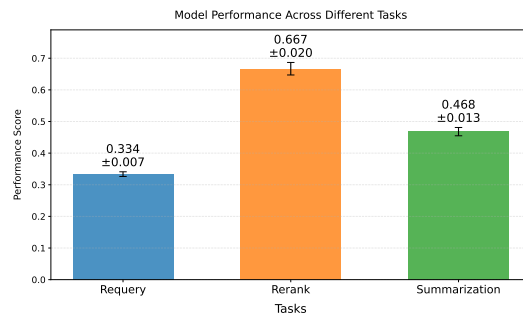
### E.1 MULTIPLE REQUERY SETTINGS

There is another evaluation setting where we could better exploit LMM’s capability. We conducted an additional experiment with LLaVA-OneVision (7B), where we generated 5 requery attempts per question (temperature = 1.0) and selected the best requery based on similarity to human-annotated queries. The result is provided in Table 5. The result showcases that conducting requery for five times could indeed improve the end-to-end performance. This provides a valuable insight for future development of a multimodal AI search engine pipeline.

**Table 5: Multiple requery evaluation setting.** We report the end-to-end performance in the table.

Model	Req. Times	Avg	News	Know.
LLaVA-OV-7B	1	29.6	33.1	19.7
LLaVA-OV-7B	5	30.1	33.3	21.3

### E.2 ROBUSTNESS EXPERIMENT AND HUMAN CORRELATION OF EVALUATION METRICS



**Figure 9: Robustness Analysis of LLaVA-OneVision-7B Performance.**

We conduct the robustness analysis comprising iterative experiments (n=8) utilizing LLaVA-OneVision-7B with a temperature parameter of 0.6 across three primary tasks: requery, rerank, and summarization. The results, visualized in Fig. 9, demonstrate remarkable stability across all tasks, with minimal variance in performance metrics, suggesting the robustness of the evaluation results.

We also evaluate the correlation between the human evaluation and the automatic evaluation metrics of the requery and summarization tasks. Three independent annotators assess the quality of the requery and summarization outputs using a score of 1, 2, 3, 4. Then the score is normalized to a [0,1] range. We compute two correlation coefficients: Pearson’s correlation coefficient ( $r$ ) and Spearman’s rank correlation coefficient ( $\rho$ ) for the significance testing. The results are shown in Table 6. The analysis revealed substantial correlations between human and automatic evaluations across both tasks. All correlations were highly significant ( $p < 0.001$ ), providing strong evidence for the validity of our automatic evaluation metrics in aligning with human judgment.

**Table 6: Correlation Analysis Results**

Statistics	Pearson’s $r$	Spearman’s $\rho$
Requery	0.5413	0.5503
Summarization	0.8112	0.7944

### E.3 SCALING TEST-TIME COMPUTE VS SCALING MODEL SIZE

Recent works such as OpenAI o1 (OpenAI, 2024a) and Li et al. (2024d) have highlighted the critical role of scaling test-time computation in enhancing model performance. Our end-to-end task, which



requires multiple Internet interactions, presents an opportunity to investigate the potential of scaling test-time computation compared to scaling model size. To explore this, we conduct experiments using LLaVA-OneVision-7B (Li et al., 2024b), focusing on scaling test-time computation, and compare against LLaVA-OneVision-72B scaling in model size, which aims to provide insights into the relative benefits of increased inference computation versus increased model parameters.

For scaling up the test-time computation, we adopt a multi-modal search strategy similar to best-of-N solution, where ‘N’ denotes 25 in our settings. Specifically, for LLaVA-OneVision-7B, we first prompt the model to generate a requery 5 times, from which we selected the one with the highest requery score  $S_{req}$ . This requery is then used to retrieve brief results from 8 websites from a search engine. The model is again prompted 5 times to select the most informative website. After removing duplicates from the selected websites, we extract the full website content from the remaining ones and prompt the model to answer 5 times, obtaining 25 end-to-end outputs in total. We compute the F1 score for each answer against the ground truth and take the maximum as the model’s end-to-end score for the query. Table 7 shows that LLaVA-OneVision-7B (TTC) achieves the score of 55.2% in the end-to-end task, significantly enhancing the original score of 29.6%, which surpasses LLaVA-OneVision-72B’s 44.9% and GPT-4V’s 52.1%. This result reveals the substantial potential of scaling test-time computation, validating the effectiveness of this technique as introduced by OpenAI o1. Our findings provide valuable insights for future research in this domain, suggesting that increased inference computation may offer comparable or superior performance improvements to increased model size not only in math and code tasks, but also in multimodal search tasks.

Table 7: **Scaling Test-Time Compute vs Scaling Model Size.** ‘TTC’ and  $S_{e2e}$  denote Test-Time Computation and the score of end-to-end task.

Model	Inference Cost	$S_{e2e}$
LLaVA-OV-7B	1	29.6
LLaVA-OV-7B (TTC)	~25	55.2
LLaVA-OV-72B	~6	44.9

#### E.4 DEFINITION OF ERRORS IN THE REQUERY AND SUMMARIZATION TASKS

Five types of requery error:

- *Lacking Specificity*, where the model fails to include all the specific information in the requery and therefore leads to sub-optimal search results. For example, the query is asking the release date of Vision Pro in China. However, the model omits the condition of China and directly asks about the release date of Vision Pro.
- *Inefficient Query*, where the model does not consider the real scenario and the requery is inefficient for the search engine to find the answer. For example, the query is asking whether the Van Gogh’s Sunflowers and Antoni Clavé’s Grand Collage are both oil paintings. Clearly, it is a commonsense that Van Gogh’s Sunflowers is an oil painting and Antoni Clavé’s Grand Collage is much less well-known. An efficient query should be asking about the images of Antoni Clavé’s Grand Collage and further determine if it is also an oil painting by directly looking at it. However, the model directly asks the original query to the search engine. There is very little chance that an exact same question has ever been raised so probably this requery will bring very little helpful information.
- *Excluding Image Search Results*, where the model totally ignores the information in the screenshot of the image search results and therefore lacks important specific information in the requery. For example, the query is ‘When did this football player obtain the gold medal?’ and provides an image of the player. The model is supposed to find out the player’s name by viewing the image search result and raise a requery like ‘[PLAYER NAME] obtained the gold medal time’. However, the model fails to incorporate the player’s name in the requery and definitely the retrieved websites will not include any helpful information.
- *No Change*, where the model just uses the question as the query input to the search engine.
- *Irrelevant*, where the model either matches wrong information from the image search result or mistakenly understands the query and outputs an irrelevant requery.

Five types of summarization error:

- *Text Reasoning Error*, where the model fails to extract the answer from the website textual information.

- *Image-text Aggregation Error*, where obtaining the answer needs combining the information from both images and texts. The model fails to do so.
- *Image reasoning Error*, where the model fails to extract the answer from the image, and the answer can only be obtained from the image.
- *Hallucination* (Huang et al., 2023), where the model provides an unfaithful answer that cannot be grounded in the given content.
- *Informal*, the output format does not follow the prompt specifications, the same error type in the end-to-end task.

## F ADDITIONAL EXPERIMENTAL DETAILS

**Rationality of the weight in final score.** Considering the sequential nature of the tasks, the reason for our weighting scheme includes two complementary perspectives:

1. The importance of each task due to the cascaded nature is already reflected in the end-to-end score. Detailedly, although it is easy to discern that the upstreamed task is more important, it is difficult to assign a precise weight to each of them. So we do not manually assign the weights but directly focus mainly on the end-to-end score, which implicitly considers their cascaded nature.
2. The individual task weights serve as complementary metrics rather than indicators of relative importance. Relying solely on end-to-end evaluation, while comprehensive, may obscure the performance characteristics of individual components and hinder targeted improvements. We therefore maintain independent evaluation of each task, with the weight distribution designed to balance the prominence of the end-to-end metric against the component-level assessments. This dual evaluation strategy enables both system-level optimization and component-specific refinements.

**More Implementation Details** All our experiments of open-source models are conducted without any fine-tuning on search data or tasks. As for the prompts, the requery prompt contains 3 examples to better guide LMMs to output a valid requery. While prompts for other tasks are all in a zero-shot setting. We prompt the LMM to output as few words as possible for a better match with the ground truth. We employ the metric introduced in Section 2.3. Besides, we recruit eight qualified college students and ask them to solve the problems in MMSEARCH independently, following the same pipeline of MMSEARCH-ENGINE. This score serves as a baseline for human performance. We conduct all experiments on NVIDIA A100 GPUs.

The input image dimensions for the webpage’s top section screenshot were set to  $1024 \times 1024$  pixels. For the full-page screenshot, we set the initial webpage width to 512 pixels, although the actual width of a small portion of webpages may vary due to its layout settings. Furthermore, considering that a full-page screenshot can be extremely lengthy, directly inputting it as a single image into an LLM would result in excessive downsizing, making the content too vague for accurate identification. To address this, we segmented the full-page screenshot into multiple images, starting from the top, with each segment measuring 512 pixels in height. Because of the context length limitations of LMMs, the maximum number of full-page screenshot segments is therefore restricted to 10.

**Full-page Screenshot Slimming.** For the full-page screenshot, we compute the Sobel gradients (Kanopoulos et al., 1988) to detect the edges and generate a gradient magnitude image. We iteratively remove the areas with gradients below a threshold, which represent the blank areas. This approach, shown in Fig. 10, effectively reduces image size while maintaining the document content.

**Model Sources.** For different LMMs, we select their latest models with size around 7B for evaluation to fully reveal their multimodal search proficiency. Table 8 presents the release time and model sources of LMMs used in MMSEARCH.

**Input Prompts of LMM for Response Generation.** We showcase the input prompts of LMM for the three tasks respectively in Table 9-11. We adopt two types of prompts for queries with an image and without images. For query with an image, we specifically require the LMM to leverage the image search result to solve the task.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

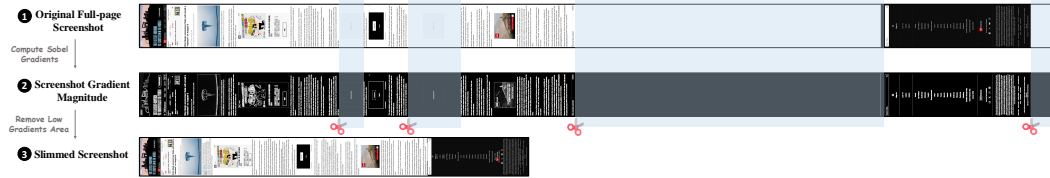


Figure 10: **Illustration of the Screenshot Slim Process.** We leverage Sobel gradients (Kanopoulos et al., 1988) to identify blank areas and remove them. After slimming, the screenshot size is largely reduced without any information loss.

Table 8: **The Release Time and Model Source of LMMs Used in MMSEARCH.**

Model	Release Time	Source
GPT-4V (OpenAI, 2023c)	2023-09	<a href="https://platform.openai.com/">https://platform.openai.com/</a>
GPT-4o (OpenAI, 2024b)	2024-05	<a href="https://platform.openai.com/">https://platform.openai.com/</a>
Claude 3.5 Sonnet (Anthropic, 2024)	2024-06	<a href="https://www.anthropic.com/news/claude-3-5-sonnet">https://www.anthropic.com/news/claude-3-5-sonnet</a>
InternLM-XC2.5 (Zhang et al., 2024a)	2024-07	<a href="https://github.com/InternLM/InternLM-XComposer">https://github.com/InternLM/InternLM-XComposer</a>
Mantis (Jiang et al., 2024b)	2024-05	<a href="https://tiger-ai-lab.github.io/Mantis/">https://tiger-ai-lab.github.io/Mantis/</a>
LLaVA-NeXT-Interleave (Li et al., 2024c)	2024-06	<a href="https://github.com/LLaVA-VL/LLaVA-NeXT">https://github.com/LLaVA-VL/LLaVA-NeXT</a>
InternVL2 (Chen et al., 2024d)	2024-07	<a href="https://github.com/OpenGVLab/InternVL">https://github.com/OpenGVLab/InternVL</a>
mPlug-Owl3 (Ye et al., 2024)	2024-08	<a href="https://github.com/X-PLUG/mPLUG-Owl">https://github.com/X-PLUG/mPLUG-Owl</a>
Idefics3 (Laurençon et al., 2024)	2024-08	<a href="https://huggingface.co/HuggingFaceM4/Idefics3-8B-Llama3">https://huggingface.co/HuggingFaceM4/Idefics3-8B-Llama3</a>
LLaVA-OneVision (Li et al., 2024b)	2024-08	<a href="https://llava-vl.github.io/blog/2024-08-05-llava-onevision/">https://llava-vl.github.io/blog/2024-08-05-llava-onevision/</a>
Qwen2-VL (Qwen Team, 2024)	2024-08	<a href="https://github.com/QwenLM/Qwen2-VL">https://github.com/QwenLM/Qwen2-VL</a>

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

Table 9: **Input Prompt of LMMs for Requery.** We adopt two different prompts for the query with image input and without image input. We leverage a 3-shot prompt to guide the LMM to generate a reasonable requery.

Question	Prompt
Query without image	<p>You are a helpful assistant. I am giving you a question, which cannot be solved without external knowledge. Assume you have access to a text-only search engine (e.g., google). Please raise a query to the search engine to search for what is useful for you to answer the question correctly. Your query needs to consider the attribute of the query to search engine. Here are 3 examples:</p> <p>Question: Did Zheng Xiuwen wear a knee pad in the women’s singles tennis final in 2024 Paris Olympics? Query to the search engine: Images of Zheng Xiuwen in the women’s singles tennis final in 2024 Paris Olympics</p> <p>Question: When will Apple release iPhone16? Query to the search engine: iPhone 16 release date</p> <p>Question: Who will sing a French song at the Olympic Games closing ceremony? Query to the search engine: Singers at the Olympic Games closing ceremony, French song.</p> <p>Question: <math>\{question\}</math>.</p> <p>Query to the search engine (do not involve any explanation):</p>
Query with image	<p>You are a helpful assistant. I am giving you a question including an image, which cannot be solved without external knowledge. Assume you have access to a search engine (e.g., google). Please raise a query to the search engine to search for what is useful for you to answer the question correctly. You need to consider the characteristics of asking questions to search engines when formulating your questions. You are also provided with the search result of the image in the question. You should leverage the image search result to raise the text query. Here are 3 examples:</p> <p>Question: Did Zheng Xiuwen wear a knee pad in the women’s singles tennis final in 2024 Paris Olympics? Query to the search engine: Images of Zheng Xiuwen in the women’s singles tennis final in 2024 Paris Olympics</p> <p>Question: When will Apple release iPhone16? Query to the search engine: iPhone 16 release date</p> <p>Question: Who will sing a French song at the Olympic Games closing ceremony? Query to the search engine: Singers at the Olympic Games closing ceremony, French song</p> <p>Question: <math>\{query\_image\}\{question\}</math>. The image search result is: <math>\{image\_search\_result\}</math></p> <p>Query to the search engine (do not involve any explanation):</p>

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

Table 10: **Input Prompt of LMMs for Rerank.** We adopt two different prompts for the query with image input and without image input.

Question	Prompt
Query without image	<p>You are a helpful assistant. I am giving you a question and 8 website information related to the question (including the screenshot, snippet and title). You should now read the screenshots, snippets and titles. Select 1 website that is the most helpful for you to answer the question. Once you select it, the detailed content of them will be provided to help you correctly answer the question. The question is <math>\{question\}</math>. The website informations is <math>\{website\_information\}</math>. You should directly output 1 website’s index that can help you most, and enclose the website in angle brackets. The output format should be: &lt;Website Index &gt;. An example of the output is: &lt;Website 1 &gt;. Your answer:</p>
Query with image	<p>You are a helpful assistant. I am giving you a question including an image. You are provided with the search result of the image in the question. And you are provided with 8 website information related to the question (including the screenshot, snippet, and title). You should now read the screenshots, snippets and titles of these websites. Select 1 website that is the most helpful for you to answer the question. Once you select it, the detailed content of them will be provided to help you correctly answer the question. The question is <math>\{query\_image\}\{question\}</math>. The image search result is <math>\{image\_search\_result\}</math>. The website information is <math>\{website\_information\}</math>. You should directly output 1 website’s index that can help you most, and enclose the website in angle brackets. The output format should be: &lt;Website Index &gt;. An example of the output is: &lt;Website 1 &gt;. Your answer:</p>



1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

Table 11: **Input Prompt of LMMs for Summarization.** We adopt two different prompts for the query with image input and without image input.

Question	Prompt
Query without image	<p>You are a helpful assistant. I am giving you a question and 1 website information related to the question. Please follow these guidelines when formulating your answer: 1. If the question contains a false premise or assumption, answer “invalid question”. 2. When answering questions about dates, use the yyyy-mm-dd format. 3. Answer the question with as few words as you can.</p> <p>You should now read the information of the website and answer the question. The website information is <math>\{website\_information\}</math>. The question is <math>\{question\}</math>. Please directly output the answer without any explanation:</p>
Query with image	<p>You are a helpful assistant. I am giving you a question including an image. You are provided with the search result of the image in the question. And you are provided with 1 website information related to the question. Please follow these guidelines when formulating your answer: 1. If the question contains a false premise or assumption, answer “invalid question”. 2. When answering questions about dates, use the yyyy-mm-dd format. 3. Answer the question with as few words as you can.</p> <p>You should now read the information of the website and answer the question. The website information is <math>\{website\_information\}</math>. The image search result is <math>\{image\_search\_result\}</math>. The question is <math>\{query\_image\}\{question\}</math>. Please directly output the answer without any explanation:</p>

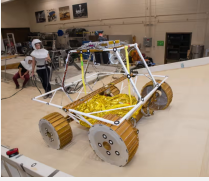
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

G MORE DATA DETAILS

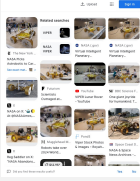
G.1 DATA EXAMPLE OF 4 EVALUATION TASKS

**Query Information**

**Image**



**Image Search Result**



**Question:**  
In which day of 2024 was the lunar rover project in the picture announced to be canceled?

**Task1 Requery**

Query Information

**LMM Requery:**  
VIPER was announced to be cancelled

**Requery Annotation:**  
NASA Ends VIPER Project

**Task2 Rerank**

Query Information

Brief Result

**<Website 1>:**  
**Title:** NASA Ends VIPER Project, Continues Moon Exploration  
**Snippet:** NASA Ends VIPER Project, Continues Moon Exploration. Tiernan P. Doyle. Jul 17, 2024. RELEASE 24-094. ... (Volatiles Investigating Polar Exploration Rover) project. NASA stated cost increases, delays to the launch date, and the risks of future cost growth as the reasons to stand down on the mission. The rover was originally planned to launch in ...

**<Website 2>:**  
**Title:** NASA cancels VIPER lunar rover - SpaceNews  
**Snippet:** The VIPER lunar rover. Credit: NASA. BUSAN, South Korea '2014 NASA has canceled a robotic lunar rover mission that would have searched for ice at the south pole of the moon, citing development ...

**<Website 3>:**  
**Title:** NASA Ends VIPER Project, Continues Moon Exploration - Yahoo Finance  
**Snippet:** NASA Ends VIPER Project, Continues Moon Exploration. PR Newswire. Wed, Jul 17, 2024, 4:13 PM 3 min read. Link Copied. 0. WASHINGTON, July 17, 2024 /PRNewswire/ -- Following a comprehensive ...

**<Website 4>:**  
**Title:** NASA axes robotic lunar rover project VIPER due to rising costs - MSN  
**Snippet:** NASA has ended its VIPER project, which was hoping to launch the agency's first robotic lunar rover to the moon, due to the increasing costs of the program. "Decisions, of course, like this are ...

**LMM Rerank:**  
<Website 2>


**Rerank Annotation:**  
**Valid:** [< Website 1>, <Website 2>, <Website 3>]  
**Unsure:** [<Website 4>, <Website 5>]  
**Invalid:** [<Website 6>, <Website 7>, <Website 8>]

**Task3 Summarization**

Query Information

Full Website Content

**Full-page Screenshot**



**Content:**  
would have searched for ice at the south pole of the moon , citing development delays and cost overruns . NASA announced July 17 that it would end development of the Volatiles Investigating Polar Exploration Rover ( VIPER ) mission . The rover , to be sent to the south polar region of moon on a commercial lander called Griffin from Astrobotic Technology , would have explored terrain including permanently shadowed regions to better understand the extent and form of water ice there . At a briefing to announce the cancellation , agency officials said costs of VIPER had grown .  
Posted inCivil NASA cancels VIPER lunar rover by Jeff Foust July 17 , 2024July 17 , 2024 Click to share on X ( Opens in new window ) Click to share on Facebook ( Opens in new window ) Click to share on LinkedIn ( Opens in new window ) Click to share on Reddit ( Opens in new window ) Click to email a link to a friend ( Opens in new window ) Click to share on Clipboard ( Opens in new window ) BUSAN , South Korea — NASA has canceled a robotic lunar rover mission that 's lander that would deliver the rover to the moon under a CLPS task order worth \$ 322 million . NASA said Griffin was now expected to be ready for the mission no earlier than September 2025 . With VIPER canceled , NASA will retain the task order for Griffin . The mission will instead become a technology demonstrator , carrying a mass simulator in place of the rover to test Griffin ' s ability to land large payloads . Kearns said NASA considered flying science payloads instead , but since the lander was designed for carrying a rover by more than 30 % , triggering a termination review by the agency . NASA had confirmed VIPER in 2021 at a cost of \$ 433.5 million . Joel Kearns , deputy associate administrator for exploration in NASA ' s Science Mission Directorate , said the latest estimate was \$ 609.6 ...

**LMM Answer:**  
July 17

**Answer Annotation:**  
07-17

**Task4 End-to-end**

Query Information

**LMM Answer:**  
Sep 12

**Answer Annotation:**  
07-17

Figure 11: Example Input, LMM Output, and Ground Truth for Four Evaluation Tasks. The color-coding of each module corresponds to Fig. 4. Task1 Requery (green), Task2 Rerank (purple), Task3 Summarization (blue), and Task4 End-to-end (yellow) are shown. Image best viewed in color.

1458 G.2 SUBFIELD DEFINITION  
1459

1460 **News** area encompasses a vast spectrum of information, ranging from everyday events to engag-  
1461 ing entertainment content and specialized fields such as scientific discoveries and financial analysis.  
1462 This comprehensive coverage serves as a rigorous assessment of the model’s ability to process in-  
1463 formation in diverse domains. We divide this expansive area into eight distinct subfields:

- 1464 • **Traditional Sports:** Data concerning traditional athletic competitions, team performances,  
1465 player statistics, and sporting events. This includes scores, league standings, player trans-  
1466 fers, and analysis of various professional sports across different leagues and countries.
- 1467 • **e-Sports:** Information about competitive video gaming, including tournament results,  
1468 player rankings, and league information. This covers various game titles, team formations,  
1469 streaming viewership statistics, and tournament information.
- 1470 • **Technology:** Information about technological innovations, gadgets, software develop-  
1471 ments, and tech industry news. This includes product launches, software updates, cyberse-  
1472 curity issues, and artificial intelligence advancements.
- 1473 • **Paper:** Content related to academic papers, research publications, and scholarly articles in  
1474 various artificial intelligence fields. The queries include method explanation, figure under-  
1475 standing, and experiment settings.
- 1476 • **Entertainment:** Data about movies, music, television, celebrities, and other forms of pop-  
1477 ular entertainment. It also includes data concerning video games.
- 1478 • **Finance:** Information on financial markets, economic indicators, business news, and mon-  
1479 etary policies. This covers stock prices, company earnings reports, company financial state-  
1480 ments, and regulatory news regarding finance.
- 1481 • **General News:** Broad coverage of various news topics not specific to any particular sub-  
1482 field. This includes a mix of local and global events, human interest stories, lifestyle arti-  
1483 cles, climate news, and general interest content that doesn’t fit neatly into other specialized  
1484 news subfields.
- 1485 • **False Premise:** Data related to misinformation or incorrect assumptions in the query. This  
1486 subfield focuses on fact-checking capabilities. All the answers to the queries of this subfield  
1487 are ‘invalid question’.

1488  
1489  
1490 **Knowledge** area represents broad subfields of information and data related to general knowledge  
1491 across various disciplines. This area concentrates on rare knowledge that most LMMs fail to answer.  
1492 We categorize this area into five subfields:

- 1493 • **Architecture:** Information about building design, architectural styles, building informa-  
1494 tion, and construction projects. This includes city landmarks, the comparison of architec-  
1495 tural styles, and multi-view architecture matchings.
- 1496 • **Arts:** Data concerning visual arts, drawings, sculptures, badges, and other forms of creative  
1497 expression. This covers artwork details, artist profiles, artwork history, and artwork style  
1498 comparisons.
- 1499 • **Fashion:** Content related to clothing trends, fashion brands, and designer collections. This  
1500 includes retail price, clothing style, release date, and brand information.
- 1501 • **Astronomy:** Information about celestial objects, space exploration, astronomical phenom-  
1502 ena, and related research. This covers observational data from telescopes and image results  
1503 from space missions. The questions focus on the background information of these celestial  
1504 objects presented in the query image.
- 1505 • **Anime:** Data about Japanese animation, including series storylines and character informa-  
1506 tion. This encompasses character background, character appearance, voice actor informa-  
1507 tion, and chapter information.
- 1508 • **Auto:** Content related to automobiles, including vehicle specifications, industry trends, and  
1509 automotive technology. This covers new car models, performance test results, coefficients  
1510 of cars, and release date.

H QUALITATIVE EXAMPLES

**Question:** In the LPL 2024 Summer season, how many teams were in Group Ascend?

**Answer:** 9

**Area:** News

**Subfield:** ⚽ Sports

Round1 Requery

**GPT-4o:** LPL 2024 Summer season Group Ascend teams count

**Qwen2-VL:** LPL 2024 Summer season Group Ascend teams

**LLaVA-OneVision:** LPL 2024 Summer season Group Ascend teams

Round2 Rerank

**GPT-4o Brief Results:**



**<Website 1>:**  
**Title:** LPL 2024 Summer - Leaguepedia | League of Legends Esports Wiki - Fandom  
**Snippet:** The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure, including introduction of "Fearless Draft" rules and return of group stages. Seventeen teams play against each other first in four double round robins, and then in two single round robins.



**<Website 5>:**  
**Title:** Group Stage / LPL 2024 Summer - schedule, results | u2014 Escorenews  
**Snippet:** LPL 2024 Summer Match results, calendar, VODs, stream, team rosters, schedule ... VODs, stream, team rosters, schedule. LoL News Bets and predictions Matches Events Teams Players. LPL 2024 Summer LoL. u2022 2024-06-01 - 0000-00-00 u2022 \$578600. Stats. Event Placements Qualifier Group Stage Playoff Regional Finals. Group Stage Playoff - Phase 3 ...



**<Website 2>:**  
**Title:** LPL Summer 2024 - Likipedia League of Legends Wiki  
**Snippet:** The LPL Summer 2024 split is the second split of the 2024 LPL season. The league maintains all 17 teams from the Spring Split, and will be held across China, in cities such as Shanghai, Suzhou, Shenzhen, Xi'an and Beijing. Bilibili Gaming is the defending title champion. This split, the LPL will experiment with a different format, featuring a ...



**<Website 6>:**  
**Title:** 2024 LPL season - Wikipedia  
**Snippet:** The 2024 LPL season is the 12th and ongoing season of the League of Legends Pro League ... The bottom two ascend group teams and the top four nirvana group teams will have to contest an additional match. [6] As per usual, the champion for Summer 2024 will qualify for the 2024 World Championship as China's number one seed. Spring. Regular Season



**<Website 3>:**  
**Title:** LPL 2024 Summer Placements - Leaguepedia | League of Legends Esports Wiki  
**Snippet:** The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure, including introduction of "Fearless Draft" rules and return of group stages. Seventeen teams play against each other first in four double round robins, and then in two single round robins.



**<Website 7>:**  
**Title:** JD Gaming vs. EDward Gaming / LPL 2024 Summer Placements - Reddit  
**Snippet:** JD Gaming vs. EDward Gaming / LPL 2024 Summer Placements - Week 4 - Group A / Post-Match Discussion LPL 2024 SUMMER ... With this win by JDG, FPX advance to Group Ascend alongside JDG for the LPL 2024 Summer Season. JDG | Leaguepedia | Likipedia | Website | Twitter EDG | Leaguepedia | Likipedia ... Team WE vs. Bilibili Gaming / LPL 2024 ...



**<Website 4>:**  
**Title:** LPL 2024 Summer LoL Coverage | GosuGamers  
**Snippet:** The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure. Placements tournament will determine the groupings for the group stage. Jun 2024. 01.



**<Website 8>:**  
**Title:** League of Legends LPL 2024 Summer Split - Sportskeeda  
**Snippet:** The LPL 2024 Summer Split format is vastly different, and a big change for a major region. The first stage of the Summer Split will be the Placements stage. Teams are divided into four groups with ...








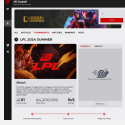
**GPT-4o Rerank:** <Website 6>

Figure 12: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619








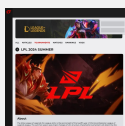
Round2 Rerank

**Qwen2-VL Brief Results:**

	<p><b>&lt;Website 1&gt;:</b>  <b>Title:</b> LPL 2024 Summer - Leaguepedia   League of Legends Esports Wiki - Fandom  <b>Snippet:</b> The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure, including introduction of "Fearless Draft" rules and return of group stages. Seventeen teams play against each other first in four double round robins, and ...</p>		<p><b>&lt;Website 5&gt;:</b>  <b>Title:</b> LPL Summer Season 2024 stats - Games of Legends  <b>Snippet:</b> View all the stats for LPL Summer Season 2024: matches result, team ranking, best players, most played champions.</p>
	<p><b>&lt;Website 2&gt;:</b>  <b>Title:</b> LPL Summer 2024 - Liquipedia League of Legends Wiki  <b>Snippet:</b> The LPL Summer 2024 is the second split of the League of Legends Pro League season. Stay up to date with match results, schedules, and broadcasts here!</p>		<p><b>&lt;Website 6&gt;:</b>  <b>Title:</b> LPL Summer 2024 - Group Stage Statistics - Liquipedia  <b>Snippet:</b> Liquipedia app major update: Revamped player and team pages with schedule, results, stats, achievements and more. Download the latest version on iOS or Android and read our update blog here.</p>
	<p><b>&lt;Website 3&gt;:</b>  <b>Title:</b> LPL/2024 Season/Summer Placements - League of Legends Esports Wiki  <b>Snippet:</b> The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of Chinas professional League of Legends league. The summer season sees major changes of LPL structure, including introduction of "Fearless Draft" rules and return of group stages. Seventeen teams play against each other first in four double round robins, and then in two single round robins. The top ten teams ...</p>		<p><b>&lt;Website 7&gt;:</b>  <b>Title:</b> LPL/2024 Season/Summer Placements/Champion Statistics  <b>Snippet:</b> Tournament: LPL/2024 Season/Summer Placements - Showing Values Per Game - Open As Query. Champion Statistics - 130 Total Games Played with 92 Champions Contested. Restrict By Role (Query Permalinks): Top.</p>
	<p><b>&lt;Website 4&gt;:</b>  <b>Title:</b> LPL 2024 Summer - Group Stage - Liquipedia League of Legends Wiki  <b>Snippet:</b> The Group Stage Draw Show will be held on May 21, 2024 20:00 CST. [1] Teams are split into 4 seed pools to ensure quality groups. Pool 4 will consist of the bottom 4 teams, namely: RNG, RA, EDG, UP. One team from Pool 1, 2 and 4 will be seeded to each group. One group will comprise of two Pool 3 teams, while the 3 other groups will have one ...</p>		<p><b>&lt;Website 8&gt;:</b>  <b>Title:</b> LPL 2024 Summer LoL Coverage   GosuGamers  <b>Snippet:</b> The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure. Placements tournament will determine the groupings for the group stage. Jun 2024. 01.</p>

**Qwen2-VL Rerank: <Website 1>**

**LLaVA-OneVision Brief Results:**

	<p><b>&lt;Website 1&gt;:</b>  <b>Title:</b> LPL Summer 2024 - Liquipedia League of Legends Wiki  <b>Snippet:</b> The LPL Summer 2024 split is the second split of the 2024 LPL season. The league maintains all 17 teams from the Spring Split, and will be held across China, in cities such as Shanghai, Suzhou, Shenzhen, Xi'an and Beijing. Bilibili Gaming is the defending title champion. This split, the LPL will experiment with a different format, featuring a ...</p>		<p><b>&lt;Website 5&gt;:</b>  <b>Title:</b> LPL Summer Season 2024 stats - Games of Legends  <b>Snippet:</b> View all the stats for LPL Summer Season 2024: matches result, team ranking, best players, most played champions. LPL Summer Season 2024 stats: Team ranking, Top KDA, picks and bans. Sign In (Log In) ... LPL Summer Season 2024 stats: Tournament data: Number of games: 152. Average game duration: 32:57. Average kills / game: 28. Shortest game:</p>
	<p><b>&lt;Website 2&gt;:</b>  <b>Title:</b> LPL 2024 Summer - Leaguepedia   League of Legends Esports Wiki  <b>Snippet:</b> The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of Chinas professional League of Legends league. The summer season sees major changes of LPL structure, including introduction of "Fearless Draft" rules and return of group stages. Seventeen teams play against each other first in four double round robins, and then in two single round robins.</p>		<p><b>&lt;Website 6&gt;:</b>  <b>Title:</b> Group Stage / LPL 2024 Summer - schedule, results - Escorenews  <b>Snippet:</b> LPL 2024 Summer Match results, calendar, VODs, stream, team rosters, schedule ... VODs, stream, team rosters, schedule. LoL News Bets and predictions Matches Events Teams Players. LPL 2024 Summer LoL. u2022 2024-06-01 - 0000-00-00 u2022 S578600. Stats. Event Placements Qualifier Group Stage Playoff Regional Finals. Group Stage Playoff - Phase 3 ...</p>
	<p><b>&lt;Website 3&gt;:</b>  <b>Title:</b> LPL 2024 Summer - Group Stage - Liquipedia League of Legends Wiki  <b>Snippet:</b> The Group Stage Draw Show will be held on May 21, 2024 20:00 CST. [1] Teams are split into 4 seed pools to ensure quality groups. Pool 4 will consist of the bottom 4 teams, namely: RNG, RA, EDG, UP. One team from Pool 1, 2 and 4 will be seeded to each group. One group will comprise of two Pool 3 teams, while the 3 other groups will have one ...</p>		<p><b>&lt;Website 7&gt;:</b>  <b>Title:</b> 2024 LPL season - Wikipedia  <b>Snippet:</b> The 2024 LPL season is the 12th and ongoing season of the League of Legends Pro League ... The bottom two ascend group teams and the top four nirvana group teams will have to contest an additional match. [6] As per usual, the champion for Summer 2024 will qualify for the 2024 World Championship as China's number one seed. Spring. Regular Season</p>
	<p><b>&lt;Website 4&gt;:</b>  <b>Title:</b> LPL 2024 Summer - Rumble Stage - Liquipedia League of Legends Wiki  <b>Snippet:</b> Format. Rumble Stage. Single round robin within each group. All matches are Bo3 using Standard Draft. For the Group "Ascend" (High Group). Top 2 teams get seeded into the semifinals. 3rd and 4th place teams get seeded into the fourth round of playoffs. 5th and 6th place teams get seeded into the third round of playoffs.</p>		<p><b>&lt;Website 8&gt;:</b>  <b>Title:</b> LPL 2024 Summer LoL Coverage   GosuGamers  <b>Snippet:</b> The 2024 League of Legends Pro League (LPL) is the second split of the twelfth year of China's professional League of Legends league. The summer season sees major changes of LPL structure. Placements tournament will determine the groupings for the group stage. Jun 2024. 01.</p>

**LLaVA-OneVision Rerank: <Website 4>**

Figure 13: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.



1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

Round3 Summarization

**GPT-4o**

**Full-page Screenshot:**



**Content:**

for Spring 2024 will qualify for the 2024 Mid-Season Invitational 's group stage while the runner-up will qualify for the 2024 Mid-Season Invitational 's play-in stage . [ 5 ] The Summer split 's format is split into three stages . The first stage will consist of a double round-robin group stage where all seventeen ( 17 ) teams are split into three groups of 4 and a group of 5 . This first stage will implement the " Fearless Draft " , a drafting format seen in China 's LD developmental league where Champions from previous games in the top four nirvana group teams advance to the single-elimination , king-of-the hill tournament . The bottom two ascend group teams and the top four nirvana group teams will have to contest an additional match . [ 6 ] As per usual , the champion for Summer 2024 will qualify for the 2024 World Championship as China 's number one seed . Spring [ edit ] Regular Season [ edit ] Spring 2024 is the second-consecutive season where the Top seed only lost one game during the entirety of the regular season . Moreover , Bilibili Gaming is the first organization 2024 LPL season Add links From Wikipedia , the free encyclopedia Sports season 2024 LPL season League LPL Sport League of Legends Duration 22 January – 20 April ( Spring ) TBD ( Summer ) Number of teams 17 Spring Split Champions Bilibili Gaming Runners-up Top Esports Season MVP Zhuo " knight " Ding ( Bilibili Gaming ) Summer Split LPL seasons ← 2023 2025 → The 2024 LPL season is the 12th and ongoing season of the League of Legends Pro League ( LPL ) , a Chinese professional esports league for the video game League of Legends . Similar to its previous yearly splits , the 2024 LPL season will be divided into two splits : Spring and Summer . The Spring Split began on 22 January and will end on 20 April 2024 for the Grand Finals . Meanwhile , the Summer Split will begin in the latter part of 2024 following the conclusion of the Mid-Season Invitational ( MSI ) 2024 . Bilibili Gaming was crowned the 2024 Spring Champions after defeating Top Esports 3–1 in the Grand Finals rematch of the Upper Bracket Finals . This was Bilibili Gaming 's first organizational title under the name Bilibili Gaming in the LPL . [ 1 ] The team previously defeated Top Esports in the Upper Bracket Finals 3–0 and Top Esports coming outstaying MSI opponents of Genies 3–1 in the Lower Bracket Finals . Bilibili Gaming is set to qualify for MSI 2024 's Group Stage alongside the champions of the LCS , LCK and LEC , while Top Esports will be joined with the runner-up teams to qualify for the play-in tournament as per the new MSI rules . Format [ edit ] The Spring Split will have seventeen ( 17 ) competing teams during the Regular Season . Each team will participate in the single round-robin tournament where all sixteen will be played in best-of-three . Similarly to its previous seasons , the top two ( 2 ) teams will advance to the playoff bracket with the Top 3 seeded teams ascending to the four-team double-elimination tournament bracket while the remaining eight ( 8 ) teams will be competing in a single-elimination bracket for all tournament slots and the top two teams from these competitions ascending to compete the four-team double-elimination bracket . [ 2 ] The champion for Spring 2024 will qualify for the 2024 Mid-Season Invitational 's group stage while the runner-up will qualify for the 2024 Mid-Season Invitational 's play-in stage . The Summer split format is split into three stages . The first stage will consist of a double round-robin group stage where all seventeen ( 17 ) teams are split into three groups of 4 and a group of 5 . This first stage will implement the " Fearless Draft " , a drafting format seen in China 's LD developmental league where Champions from previous games in the best-of-three series are banned from the re-drafting game . The top two teams in each group and the first place team in the four-team group qualify for the " seeded " group in the second stage , while the rest qualify for the " unseeded " group . The second stage will see the teams in each of the two groups face each other in single round-robin regular draft competition in best-of-threes . The top two seeded teams in the second group advance to the four-team double-elimination tournament bracket as usual , while the remaining unseeded group teams and the top four nirvana group teams advance to the single-elimination king-of-the-hill tournament . The bottom two seeded group teams and the top four nirvana group teams will have to contest an additional match . [ 6 ] As per usual , the champion for Summer 2024 will qualify for the 2024 World Championship as China 's number one seed . [ 7 ] [ 8 ] [ 9 ] [ 10 ] [ 11 ] [ 12 ] [ 13 ] [ 14 ] [ 15 ] [ 16 ] [ 17 ] [ 18 ] [ 19 ] [ 20 ] [ 21 ] [ 22 ] [ 23 ] [ 24 ] [ 25 ] [ 26 ] [ 27 ] [ 28 ] [ 29 ] [ 30 ] [ 31 ] [ 32 ] [ 33 ] [ 34 ] [ 35 ] [ 36 ] [ 37 ] [ 38 ] [ 39 ] [ 40 ] [ 41 ] [ 42 ] [ 43 ] [ 44 ] [ 45 ] [ 46 ] [ 47 ] [ 48 ] [ 49 ] [ 50 ] [ 51 ] [ 52 ] [ 53 ] [ 54 ] [ 55 ] [ 56 ] [ 57 ] [ 58 ] [ 59 ] [ 60 ] [ 61 ] [ 62 ] [ 63 ] [ 64 ] [ 65 ] [ 66 ] [ 67 ] [ 68 ] [ 69 ] [ 70 ] [ 71 ] [ 72 ] [ 73 ] [ 74 ] [ 75 ] [ 76 ] [ 77 ] [ 78 ] [ 79 ] [ 80 ] [ 81 ] [ 82 ] [ 83 ] [ 84 ] [ 85 ] [ 86 ] [ 87 ] [ 88 ] [ 89 ] [ 90 ] [ 91 ] [ 92 ] [ 93 ] [ 94 ] [ 95 ] [ 96 ] [ 97 ] [ 98 ] [ 99 ] [ 100 ] [ 101 ] [ 102 ] [ 103 ] [ 104 ] [ 105 ] [ 106 ] [ 107 ] [ 108 ] [ 109 ] [ 110 ] [ 111 ] [ 112 ] [ 113 ] [ 114 ] [ 115 ] [ 116 ] [ 117 ] [ 118 ] [ 119 ] [ 120 ] [ 121 ] [ 122 ] [ 123 ] [ 124 ] [ 125 ] [ 126 ] [ 127 ] [ 128 ] [ 129 ] [ 130 ] [ 131 ] [ 132 ] [ 133 ] [ 134 ] [ 135 ] [ 136 ] [ 137 ] [ 138 ] [ 139 ] [ 140 ] [ 141 ] [ 142 ] [ 143 ] [ 144 ] [ 145 ] [ 146 ] [ 147 ] [ 148 ] [ 149 ] [ 150 ] [ 151 ] [ 152 ] [ 153 ] [ 154 ] [ 155 ] [ 156 ] [ 157 ] [ 158 ] [ 159 ] [ 160 ] [ 161 ] [ 162 ] [ 163 ] [ 164 ] [ 165 ] [ 166 ] [ 167 ] [ 168 ] [ 169 ] [ 170 ] [ 171 ] [ 172 ] [ 173 ] [ 174 ] [ 175 ] [ 176 ] [ 177 ] [ 178 ] [ 179 ] [ 180 ] [ 181 ] [ 182 ] [ 183 ] [ 184 ] [ 185 ] [ 186 ] [ 187 ] [ 188 ] [ 189 ] [ 190 ] [ 191 ] [ 192 ] [ 193 ] [ 194 ] [ 195 ] [ 196 ] [ 197 ] [ 198 ] [ 199 ] [ 200 ] [ 201 ] [ 202 ] [ 203 ] [ 204 ] [ 205 ] [ 206 ] [ 207 ] [ 208 ] [ 209 ] [ 210 ] [ 211 ] [ 212 ] [ 213 ] [ 214 ] [ 215 ] [ 216 ] [ 217 ] [ 218 ] [ 219 ] [ 220 ] [ 221 ] [ 222 ] [ 223 ] [ 224 ] [ 225 ] [ 226 ] [ 227 ] [ 228 ] [ 229 ] [ 230 ] [ 231 ] [ 232 ] [ 233 ] [ 234 ] [ 235 ] [ 236 ] [ 237 ] [ 238 ] [ 239 ] [ 240 ] [ 241 ] [ 242 ] [ 243 ] [ 244 ] [ 245 ] [ 246 ] [ 247 ] [ 248 ] [ 249 ] [ 250 ] [ 251 ] [ 252 ] [ 253 ] [ 254 ] [ 255 ] [ 256 ] [ 257 ] [ 258 ] [ 259 ] [ 260 ] [ 261 ] [ 262 ] [ 263 ] [ 264 ] [ 265 ] [ 266 ] [ 267 ] [ 268 ] [ 269 ] [ 270 ] [ 271 ] [ 272 ] [ 273 ] [ 274 ] [ 275 ] [ 276 ] [ 277 ] [ 278 ] [ 279 ] [ 280 ] [ 281 ] [ 282 ] [ 283 ] [ 284 ] [ 285 ] [ 286 ] [ 287 ] [ 288 ] [ 289 ] [ 290 ] [ 291 ] [ 292 ] [ 293 ] [ 294 ] [ 295 ] [ 296 ] [ 297 ] [ 298 ] [ 299 ] [ 300 ] [ 301 ] [ 302 ] [ 303 ] [ 304 ] [ 305 ] [ 306 ] [ 307 ] [ 308 ] [ 309 ] [ 310 ] [ 311 ] [ 312 ] [ 313 ] [ 314 ] [ 315 ] [ 316 ] [ 317 ] [ 318 ] [ 319 ] [ 320 ] [ 321 ] [ 322 ] [ 323 ] [ 324 ] [ 325 ] [ 326 ] [ 327 ] [ 328 ] [ 329 ] [ 330 ] [ 331 ] [ 332 ] [ 333 ] [ 334 ] [ 335 ] [ 336 ] [ 337 ] [ 338 ] [ 339 ] [ 340 ] [ 341 ] [ 342 ] [ 343 ] [ 344 ] [ 345 ] [ 346 ] [ 347 ] [ 348 ] [ 349 ] [ 350 ] [ 351 ] [ 352 ] [ 353 ] [ 354 ] [ 355 ] [ 356 ] [ 357 ] [ 358 ] [ 359 ] [ 360 ] [ 361 ] [ 362 ] [ 363 ] [ 364 ] [ 365 ] [ 366 ] [ 367 ] [ 368 ] [ 369 ] [ 370 ] [ 371 ] [ 372 ] [ 373 ] [ 374 ] [ 375 ] [ 376 ] [ 377 ] [ 378 ] [ 379 ] [ 380 ] [ 381 ] [ 382 ] [ 383 ] [ 384 ] [ 385 ] [ 386 ] [ 387 ] [ 388 ] [ 389 ] [ 390 ] [ 391 ] [ 392 ] [ 393 ] [ 394 ] [ 395 ] [ 396 ] [ 397 ] [ 398 ] [ 399 ] [ 400 ] [ 401 ] [ 402 ] [ 403 ] [ 404 ] [ 405 ] [ 406 ] [ 407 ] [ 408 ] [ 409 ] [ 410 ] [ 411 ] [ 412 ] [ 413 ] [ 414 ] [ 415 ] [ 416 ] [ 417 ] [ 418 ] [ 419 ] [ 420 ] [ 421 ] [ 422 ] [ 423 ] [ 424 ] [ 425 ] [ 426 ] [ 427 ] [ 428 ] [ 429 ] [ 430 ] [ 431 ] [ 432 ] [ 433 ] [ 434 ] [ 435 ] [ 436 ] [ 437 ] [ 438 ] [ 439 ] [ 440 ] [ 441 ] [ 442 ] [ 443 ] [ 444 ] [ 445 ] [ 446 ] [ 447 ] [ 448 ] [ 449 ] [ 450 ] [ 451 ] [ 452 ] [ 453 ] [ 454 ] [ 455 ] [ 456 ] [ 457 ] [ 458 ] [ 459 ] [ 460 ] [ 461 ] [ 462 ] [ 463 ] [ 464 ] [ 465 ] [ 466 ] [ 467 ] [ 468 ] [ 469 ] [ 470 ] [ 471 ] [ 472 ] [ 473 ] [ 474 ] [ 475 ] [ 476 ] [ 477 ] [ 478 ] [ 479 ] [ 480 ] [ 481 ] [ 482 ] [ 483 ] [ 484 ] [ 485 ] [ 486 ] [ 487 ] [ 488 ] [ 489 ] [ 490 ] [ 491 ] [ 492 ] [ 493 ] [ 494 ] [ 495 ] [ 496 ] [ 497 ] [ 498 ] [ 499 ] [ 500 ] [ 501 ] [ 502 ] [ 503 ] [ 504 ] [ 505 ] [ 506 ] [ 507 ] [ 508 ] [ 509 ] [ 510 ] [ 511 ] [ 512 ] [ 513 ] [ 514 ] [ 515 ] [ 516 ] [ 517 ] [ 518 ] [ 519 ] [ 520 ] [ 521 ] [ 522 ] [ 523 ] [ 524 ] [ 525 ] [ 526 ] [ 527 ] [ 528 ] [ 529 ] [ 530 ] [ 531 ] [ 532 ] [ 533 ] [ 534 ] [ 535 ] [ 536 ] [ 537 ] [ 538 ] [ 539 ] [ 540 ] [ 541 ] [ 542 ] [ 543 ] [ 544 ] [ 545 ] [ 546 ] [ 547 ] [ 548 ] [ 549 ] [ 550 ] [ 551 ] [ 552 ] [ 553 ] [ 554 ] [ 555 ] [ 556 ] [ 557 ] [ 558 ] [ 559 ] [ 560 ] [ 561 ] [ 562 ] [ 563 ] [ 564 ] [ 565 ] [ 566 ] [ 567 ] [ 568 ] [ 569 ] [ 570 ] [ 571 ] [ 572 ] [ 573 ] [ 574 ] [ 575 ] [ 576 ] [ 577 ] [ 578 ] [ 579 ] [ 580 ] [ 581 ] [ 582 ] [ 583 ] [ 584 ] [ 585 ] [ 586 ] [ 587 ] [ 588 ] [ 589 ] [ 590 ] [ 591 ] [ 592 ] [ 593 ] [ 594 ] [ 595 ] [ 596 ] [ 597 ] [ 598 ] [ 599 ] [ 600 ] [ 601 ] [ 602 ] [ 603 ] [ 604 ] [ 605 ] [ 606 ] [ 607 ] [ 608 ] [ 609 ] [ 610 ] [ 611 ] [ 612 ] [ 613 ] [ 614 ] [ 615 ] [ 616 ] [ 617 ] [ 618 ] [ 619 ] [ 620 ] [ 621 ] [ 622 ] [ 623 ] [ 624 ] [ 625 ] [ 626 ] [ 627 ] [ 628 ] [ 629 ] [ 630 ] [ 631 ] [ 632 ] [ 633 ] [ 634 ] [ 635 ] [ 636 ] [ 637 ] [ 638 ] [ 639 ] [ 640 ] [ 641 ] [ 642 ] [ 643 ] [ 644 ] [ 645 ] [ 646 ] [ 647 ] [ 648 ] [ 649 ] [ 650 ] [ 651 ] [ 652 ] [ 653 ] [ 654 ] [ 655 ] [ 656 ] [ 657 ] [ 658 ] [ 659 ] [ 660 ] [ 661 ] [ 662 ] [ 663 ] [ 664 ] [ 665 ] [ 666 ] [ 667 ] [ 668 ] [ 669 ] [ 670 ] [ 671 ] [ 672 ] [ 673 ]

Figure 14: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.


31

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

Round3 Summarization

**Qwen2-VL**

**Full-page Screenshot:**



**Content:**

twelfth year of China 's professional League of Legends league . The summer season sees major changes of LPL structure , including introduction of " Fearless Draft " rules and return of group stages . Seventeen teams play against each other first in four double round robins , and then in two single round robins . The top ten teams advance to the playoffs . Overview [ ] Format [ ] Two groups , Group Ascend having 9 teams and Group Nirvana having 8 Single Round Robin within group Matches are best of three The Top 7 teams in Group

in : Chinese Tournaments , Competitions , Premier Events LPL 2024 Summer < LPL | 2024 Season Sign in to edit History Talk ( 0 ) European Pro League Season 3 vs 2 September 2024 12:00:00 +0000 LIVE • PRM 2nd Div 2025 Spring Promotion vs 2 September 2024 16:00:00 +0000 LIVE • arrMY Summer League 2024 vs 2 September 2024 16:00:00 +0000 LIVE • arrMY Summer League 2024 vs 2 September 2024 17:00:00 +0000 LIVE • arrMY Summer League 2024 to Google Calendar Social Media & Links Contents 1 Overview 1.1 Format 2 Participants 2.1 Group Ascend 2.2 Group Nirvana 3 Results 4 Match Schedule 5 VODs & Match Links 6 Individual Awards 6.1 " Man of the Match " Standings 6.2 Weekly Award 6.3 Season Awards 7 Media 7.1 Streams 7.2 Broadcast Talent 7.2.1 English 7.2.2 Mandarin 7.2.2.1 Stage / Studio Hosts 7.2.2.2 Studio Hosts 7.2.2.3 Casters 7.2.2.4 Guests Casters 7.3 Additional Content 7.4 Viewership Statistics 7.5 Announcements 8 Home Venues 9 References The 2024 League of Legends Pro League ( LPL ) is the second split of the

Ascend qualify for Playoffs The Bottom 2 teams in Group Ascend and Top 4 teams in Group Nirvana qualify for Play-in Stage The Bottom 4 teams in Group Nirvana are not qualified for playoffs Show Tiebreaker Rules Hide Tiebreaker Rules If two teams have the same number of series won , ties will be broken by : Game Differential Head to Head record Participants [ ] Show Rosters Hide Rosters Rosters By Game Player Chart Group Ascend [ ] Anyone 's Legend Ale Croco Shanks Hope Kael Tabe Qingsi Bilibili Gaming Bin Xun Wei Knight Elk Invictus Gaming YSKM glfs neny Ahn Vampire Rashomon Oh My God Hery Tianzhen Angel Starry ppgod noname Geitang Rare Atom Xiaoxu Xiaohao VicLa Assum Jwei Deceit JMZ RNG Juice Geju XBY Tangyuan Xzz huanfeng Iwandy Ming Teacherma May Team WE Wayward Yanxiang FoFo Able Mark chengz Zoom TT Gaming HOYA Beichuan ucal 1xn Feather AFei Ultra Prime Qingtian H4cker Yuekai Doggo Niket Xiaobai Yuzhang Results [ ] Group Ascend Legend Round 4 Seed Round 3 Seed Round 2 Seed Round 1 Seed Play-In Seed Team Series vs 2 September 2024 18:15:00 +0000 LIVE • arrMY Summer League 2024 vs 2 September 2024 18:15:00 +0000 LIVE • arrMY Summer League 2024 vs 2 September 2024 19:00:00 +0000 LIVE • NACL 2025 Spring Promotion vs 2 September 2024 22:00:00 +0000 LIVE • Emerald League Special Edition vs 2 September 2024 23:00:00 +0000 LIVE • NACL 2025 Spring Promotion vs 3 September 2024 00:00:00 +0000 LIVE • LCK CL 2024 Summer Playoffs vs 1 Playoffs Split 2 Replacements 2 Split 2 Playoffs Split 3 Grand Finals 2024 Season Overview Placements Split 1 Replacements 1 Split 1 Playoffs Split 2 Replacements 2 Split 2 Playoffs Split 3 Grand Finals Overview Spring Season Spring Playoffs Summer Placements Summer Season Summer Playoffs Regional Finals Championship Points Overview Team Rosters Picks & Bans Scoreboards Runes Match History Champion Stats Player Stats LPL 2024 Summer Season Tournament Information Organizer TJ Sports Location & Dates Region CN China Event Type Offline Country China Start Date 2024-07-05 End Date 2024-07-31 Broadcast Streams Twitch Tencent Full List Schedule Spoiler-Free Schedule Export • Emerald League Special Edition vs 5 September 2024 23:00:00 +0000 LIVE • NACL 2024 Spring Promotion TBD vs TBD 6 September 2024 00:00:00 +0000 LIVE • LCK CL 2024 Summer Playoffs vs TBD 6 September 2024 08:00:00 +0000 LIVE • European Pro League Season 3 vs 6 September 2024 10:00:00 +0000 LIVE • European Pro League Season 3 vs 6 September 2024 12:00:00 +0000 LIVE • PRM 2nd Div 2025 Spring Promotion TBD vs TBD 6 September 2024 16:00:00 +0000 LIVE • LFL 2025 Promotion

**Qwen2-VL Summarize: 9**

Figure 15: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

**Round3 Summarization**

### LLaVA-OneVision

**Full-page Screenshot:**



**Content:**


register and log in to edit our pages . You can also join our # leagueoflegends contributor community on our Discord . Liquipedia app match pages updated ! Liquipedia app match pages are overhauled ! Download on Android or iOS ! Liquipedia app 's match pages got completely revamped with game data , standings , VODs and more ! Download the the latest version on iOS or Android and read our update blog here . LPL 2024 Summer - Rumble Stage From Liquipedia League of Legends Wiki < LPL | 2024/Summer v Group " Ascend " ( High Group ) 2.2 Rumble Stage Group " Nirvana " ( Low Group ) 2.3 Detailed Results 2.3.1 Week 1 2.3.2 Week 2 2.3.3 Week 3 2.3.4 Week 4 2.3.5 Week 5 3 References Format [ edit ] Rumble Stage Single round robin within each group All matches are Bo3 using Standard Draft . For the Group " Ascend " ( High Group ) : Top 2 teams get seeded into the semifinals 3rd and 4th place teams get seeded into the fourth round of playoffs 5th and 6th place teams get seeded into the d e LPL & LDL/LSPL 2024 LPL Group " Ascend " ( High Group ) [ edit ] Group " Ascend " Week 5 Week 1 Week 2 Week 3 Week 4 Week 5 1 . Anyone 's Legend 1-0 2-0 +2 1 . LNG Esports 1-0 2-0 +2 1 . Weibo Gaming 1-0 2-0 +2 Bilibili Gaming 0-0 0-0 0 LGD Gaming 0-0 0-0 0 Top Esports 0-0 0-0 0 7 . FunPlus Phoenix 0-1 0-2 -2 7 . JD Gaming 0-1 0-2 -2 7 . Ninjas in Pyjamas 0-1 0-2 -2 1 . LNG Esports 3-0 6-1 +5 2 . Anyone 's Legend ▼ 1 3-1 2024 - 17:00 CST Match Page 36:33 22:46 MVPs : Bin , knight Bans Game 1 Game 2 JDG 0 2 TES JD Gaming 0 2 Top Esports July 31 , 2024 - 19:00 CST Match Page 41:08 26:12 MVPs : 369 , 369 Bans Game 1 Game 2 Match was held in Beijing . References [ edit ] Retrieved from " https : //liquipedia.net/leagueoflegends/index.php ? title=LPL/2024/Summer/Rum ble Stage & oldid=828570 " Hidden categories : Pages reading from original match table Pages storing into original game table Pages storing into original match table Do you want to help ? Just LNG LGD Gaming 1 2 LNG Esports July 27 , 2024 - 15:10 CST Match Page 32:40 30:05 30:02 MVPs : haichao , Hang , Weiwei Bans Game 1 Game 2 Game 3 Match was held in Suzhou . EDG 2 0 UP EDward Gaming 2 0 Ultra Prime July 27 , 2024 - 17:10 CST Match Page 29:19 34:56 MVPs : Cryin , Jiejie Bans Game 1 Game 2 FPX 0 2 TES FunPlus Phoenix 0 2 Top Esports July 27 , 2024 - 19:30 CST Match Page 33:43 24:02 MVPs : 19:10 CST Match Page 26:56 27:10 MVPs : Zika , Weiwei Bans Game 1 Game 2 Match was held in Shenzhen . Week 2 [ edit ] Week 2 July 8 , 2024 EDG 0 2 WE EDward Gaming 0 2 Team WE July 8 , 2024 - 17:10 CST Match Page 27:19 29:46 MVPs : Yanxiang , Able Bans Game 1 Game 2 Match was held in Xi'an . LGD 0 2 AL LGD Gaming 0 2 Anyone 's Legend July 8 , 2024 - 19:10 CST Match Page 28:09 28:36 MVPs , 2024 - 17:10 CST Match Page Spring Summer Regional Finals Championship Points LDL Seeding Stage 1 Stage 2 Stage 3 Season Finals Others Demacia Cup Legend Cup LCC Click on the " Show " link on the right to see the full list 2023 LPL Spring Summer Regional Finals Championship Points LDL Stage 1 Stage 2 Stage 3 Others Demacia Cup LCC 2022 LPL Spring Summer Regional Finals Championship Points LDL Spring Summer Others Demacia Cup LCC 2021 LPL Spring Summer Regional Finals Championship Points LDL Spring Summer Others Demacia Cup LCC 2020 LPL Spring Summer Regional Finals Championship LNG Esports July 7 , 2024 - 33:22 41:37 MVPs : Xiaohao , Xiaoxu Bans Game 1 Game 2 JDG 0 2 WBG JD Gaming 0 2 Weibo Gaming July 6 , 2024 - 19:10 CST Match Page 43:33 33:25 MVPs : Tarzan , Breathe Bans Game 1 Game 2 Match was held in Beijing . July 7 , 2024 UP 0 2 iG Ultra Prime 0 2 Invictus Gaming July 7 , 2024 - 17:10 CST Match Page 31:35 39:02 MVPs : Glfs , Ahn Bans Game 1 Game 2 NIP 0 2 LNG Ninjas in Pyjamas 0 2 July 27 , 2024 LGD 1 2

**LLaVA-OneVision Summarize: 8 ❌**

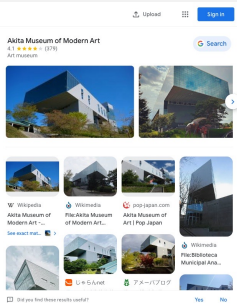
Figure 16: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

**Image:**



**Image Search Result:**



**Question:**

As a college student, how much should I pay for going to this building?

**Answer:** 210 yen

**Area:** Knowledge

**Subfield:** Architecture

**Round1 Requery**

**GPT-4o:** Akita Museum of Modern Art college student admission fee

**Qwen2-VL:** Akita Museum of Modern Art ticket price

**LLaVA-OneVision:** Akiha Museum of Modern Art

**Round2 Rerank**

**GPT-4o Brief Results:**

<p><b>&lt;Website 1&gt;</b>  <b>Title:</b> Locations, hours, and admission   MoMA  <b>Snippet:</b> Museum of Modern Art locations, hours, and admission Caption: The Museum of Modern Art Renovation and Expansion Designed by Diller Scofidio + Renfro in collaboration with Gensler. ... Students Full-time with ID, including international students \$17 Children 16 and under Free Members Free ... Admission fees for MoMA PSI are suggested. MoMA ...</p>	<p><b>&lt;Website 5&gt;</b>  <b>Title:</b> Akita Museum of Art - Akita City, Akita - Japan Travel  <b>Snippet:</b> The Akita Museum of Art is listed as the #1 thing to do in Akita by TripAdvisor. At just 310 yen for entrance (210 for college students), and with a free public gallery on the first floor, the art on display here is certainly worth a visit. This museum is also one of the few around the world that houses a large collection of work by Fujita ...</p>
<p><b>&lt;Website 2&gt;</b>  <b>Title:</b> Discounts   MoMA  <b>Snippet:</b> Discounted admission. MoMA is a participant of the following passes. For redemption instructions, please refer to the pass's website. RockMoMA. Save \$10 and enjoy two iconic NYC attractions. Start your morning with a visit to Top of the Rock, 70 floors above Rockefeller Center. After taking in the sights, make your way to MoMA and explore six floors of modern and contemporary art. New York ...</p>	<p><b>&lt;Website 6&gt;</b>  <b>Title:</b> Akita Museum of Modern Art   Places &amp; Experiences posted by Locals ...  <b>Snippet:</b> Located in Akita City, the Akita Museum of Modern Art is a wonderful museum showcasing modern art and history in the Tohoku region of Japan. It houses one of the most varied collections in the country, featuring artworks of various mediums and time periods, as well as archeological artifacts. The museum was established in 1983 and since has grown to contain over 3,000 pieces of art.</p>
<p><b>&lt;Website 3&gt;</b>  <b>Title:</b> Akita Museum of Modern Art   Yokote City Tourism Promotion Organization  <b>Snippet:</b> In 1994, the Akita Museum of Modern Art opened in the village of Akita Furusato collects and displays excellent works of art from the modern era, including the Western-style painting "Akita orchid" painted by the Akita lord and vassals in the Edo period. ... Admission free. Exhibition fee (Different depending on the content) Annual Pass ...</p>	<p><b>&lt;Website 7&gt;</b>  <b>Title:</b> MoMA Membership   Access - Museum of Modern Art  <b>Snippet:</b> Additional membership categories and discounts for artists and students are available. Call (888) 999-8861 to learn more. Memberships are not refundable or transferable. Member benefits are for personal, noncommercial use only. Categories, benefits, and prices are subject to change. Unlimited free admission to MoMA without waiting in ticket lines.</p>
<p><b>&lt;Website 4&gt;</b>  <b>Title:</b> Akita Museum of Modern Art   Museums   Japan Cultural Expo - Nihonhaku ...  <b>Snippet:</b> The Akita Museum of Modern Art opened its doors on April 20, 1994, within the Akita Furusato Village. ... Opening hours: 9:30-17:00 (last admission at 16:30) Closed Closed during the New Year holidays and maintenance (10 days in late January). Admissions Free admission. Entrance fee required for special exhibitions (varies by exhibition ...)</p>	<p><b>&lt;Website 8&gt;</b>  <b>Title:</b> MoMA Membership   Annual Pass - Museum of Modern Art  <b>Snippet:</b> Unlimited free admission to MoMA without waiting in ticket lines. Members-only gallery talks, exhibition previews, \$5 guest tickets, free films, great discounts, exclusive digital content, and more! ... Additional membership categories and discounts for artists and students are available. Call (888) 999-8861 to learn more. Memberships are not ...</p>

**GPT-4o Rerank: <Website 5>**


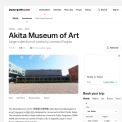






Figure 17: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

34

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889









Round2 Rerank

**Qwen2-VL Brief Results:**

 <p><b>&lt;Website 1&gt;</b>  <b>Title:</b> Akita Museum of Art - Must-See, Access, Hours &amp; Price  <b>Snippet:</b> This art museum was opened in 2013 in "Area Nakaichi" near Akita Station as a hub of arts and cultures. The modern building was designed by Tadao Ando. The interior with repeated triangles and inorganic concrete walls creates out-of-the-ordinary artistic space.</p>	 <p><b>&lt;Website 5&gt;</b>  <b>Title:</b> Akita Museum of Art - Akita Travel - japan-guide.com  <b>Snippet:</b> The Akita Museum of Art (u79cbu7530u770c u7acb u78e u853 u9928, Akita Kenritsu Bijutsukan) is an art museum in Akita City designed by renowned architect Ando Tadao. The museum exhibits a large collection of work by Fujita Tsuguharu (1886-1968), also known as Leonard Foujita, who is arguably Japan's most famous Western style painter.</p>
 <p><b>&lt;Website 2&gt;</b>  <b>Title:</b> Akita Museum of Modern Art   Yokote City Tourism Promotion Organization  <b>Snippet:</b> In 1994, the Akita Museum of Modern Art opened in the village of Akita Furusado collects and displays excellent works of art from the modern era, including the Western-style painting "Akita orchid" painted by the Akita lord and vassals in the Edo period. I will.</p>	 <p><b>&lt;Website 6&gt;</b>  <b>Title:</b> Akita Museum of Modern Art - Wikipedia  <b>Snippet:</b> Akita Museum of Modern Art (u79cbu7530u770c u7acb u8fd1 u4ec3 u78e u853 u9928, Akita Kenritsu Kindai Bijutsukan) opened in Yokote, Akita Prefecture, Japan in 1994 and houses an important collection of Akita ranga .</p>
 <p><b>&lt;Website 3&gt;</b>  <b>Title:</b> Akita Museum of Modern Art Opening hours, specific address, ticket ...  <b>Snippet:</b> How does Akita Museum of Modern Art play? Including Akita Museum of Modern Art ticket information, prices and fees, reservations, opening hours, transportation, maps, addresses, travel guides, recommended accommodation, nearby restaurants and other practical information, book Akita Museum of Modern Art tickets at super-valuable prices through Trip.com to enjoy Super discounts or learn more ...</p>	 <p><b>&lt;Website 7&gt;</b>  <b>Title:</b> Akita Museum of Modern Art   Places &amp; Experiences posted by Locals ...  <b>Snippet:</b> The Akita Museum of Modern Art is open from 10am until 5pm every day, and admission is free. Whether you're an art enthusiast or a casual visitor, the museum is an excellent option for exploring the history and cultural significance of Akita's art scene. [Text: @AI, Translation: @AI]</p>
 <p><b>&lt;Website 4&gt;</b>  <b>Title:</b> Akita Museum of Modern Art Details / Explore! Japan Travel by NAVITIME ...  <b>Snippet:</b> This museum was opened in 1994 in Yokote City, Akita Prefecture. It primarily exhibits modern works by Akita native or related artists. In addition, the Hi-Vision Gallery holds the museum's collection. The first floor has a Museum School offering regular lectures on art and many people participate.</p>	 <p><b>&lt;Website 8&gt;</b>  <b>Title:</b> Akita Museum of Art - Travel in Japan  <b>Snippet:</b> The Akita Museum of Art in Akita City, Japan, showcases modern and contemporary Japanese art, including over 5,000 pieces. Highlights include works by Tsuguharu Fujita and other renowned artists. Designed by architect Tadao Ando, the museum offers stunning exhibits and temporary exhibitions in a striking water and concrete building.</p>

**Qwen2-VL Rerank: <Website 3>**

**LLaVA-OneVision Brief Results:**

 <p><b>&lt;Website 1&gt;</b>  <b>Title:</b> Akita Museum of Modern Art - Wikipedia  <b>Snippet:</b> Akita Museum of Modern Art. Coordinates: 39 u006017 u203233 u2033N 140 u006032 u203251 u2033E. Akita Museum of Modern Art (u79cbu7530u770c u7acb u8fd1 u4ec3 u78e u853 u9928, General information. Address. Tomiigasawa-62-46, Akasaka. Town or city.</p>	 <p><b>&lt;Website 5&gt;</b>  <b>Title:</b> Akita Museum of Modern Art - All You Need to Know BEFORE ... - Tripadvisor  <b>Snippet:</b> Akita Museum of Modern Art. See all things to do. See all things to do. Akita Museum of Modern Art. 3.5. 19 reviews #12 of 55 things to do in Yokote. Art Museums. Write a review. Full view. All photos (39) Suggest edits to improve what we show. Improve this listing. Top ways to experience nearby attractions. Samurai Food Walk in Kakunodate. 2.</p>
 <p><b>&lt;Website 2&gt;</b>  <b>Title:</b> MoMA  <b>Snippet:</b> MoMA is a place that fuels creativity, ignites minds, and provides inspiration. Its extraordinary exhibitions and collection of modern and contemporary art are dedicated to helping you understand and enjoy the art of our time. Caption: The Museum of Modern Art Renovation and Expansion Designed by Diller Scofidio + Renfro in collaboration with Gensler.</p>	 <p><b>&lt;Website 6&gt;</b>  <b>Title:</b> Iraq's Modern Art Collection, Waiting to Re-emerge  <b>Snippet:</b> Such is the state of Iraq's modern art collection, renamed the National Museum of Modern Art in 2006 yet still an institution that exists mostly as an idea. That it exists at all is owed largely ...</p>
 <p><b>&lt;Website 3&gt;</b>  <b>Title:</b> The Complete Guide to the Museum of Modern Art in New York City - TripSavvy  <b>Snippet:</b> New York City's Museum of Modern Art, a.k.a. MoMA, was closed from June until October 2019 for a major \$450 million overhaul, both structurally and content-wise. Now expanded by 47,000 square feet (a 30-percent increase), MoMA boasts technological enhancements, including more than 200 digital audio guides you can listen to via free Wi-Fi, and an ever-changing temporary and permanent collection ...</p>	 <p><b>&lt;Website 7&gt;</b>  <b>Title:</b> MoMA Visitor Tips and Advice - TripSavvy  <b>Snippet:</b> The Museum of Modern Art, or MoMA for short, is one of New York City's most impressive museums. Designed by Japanese architect Yoshio Taniguchi, the building itself is as magnificent as its extensive art collection. The large six-story museum is home to many rotating exhibits and permanent art collections, so make sure to plan ahead to maximize your time.</p>
 <p><b>&lt;Website 4&gt;</b>  <b>Title:</b> Locations, hours, and admission   MoMA  <b>Snippet:</b> 10:30 a.m.-5:30 p.m. MoMA is open every day except Thanksgiving and Christmas. Exhibition and collection galleries close approximately 15 minutes prior to Museum closing. Our stores and restaurants have different hours.</p>	 <p><b>&lt;Website 8&gt;</b>  <b>Title:</b> AD Classics: The Museum of Modern Art   ArchDaily  <b>Snippet:</b> Share. The entrance to the Museum of Modern Art is tucked beneath a demure facade of granite and glass in Midtown Manhattan. Its clean, regular planes mark Yoshio Taniguchi's 2004 addition to the ...</p>

**LLaVA-OneVision Rerank: <Website 1>**

Figure 18: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

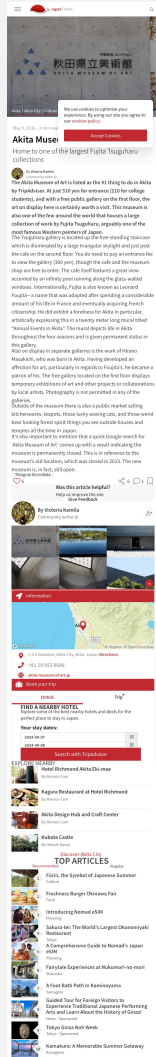


1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

Round3 Summarization

**GPT-4o**

**Full-page Screenshot:**



**Content:**

Akita City Culture Akita Akita City Culture May 9, 2016 - 2 min read Akita Museum of Art Home to one of the largest Fujita Tsugaru collections By Victoria Kamila Community writer The Akita Museum of Art is listed as the # 1 thing to do in Akita by TripAdvisor . At just 310 yen for entrance ( 210 for college students ) , and with a free public gallery on the first floor , the art on display here is certainly worth a visit . This museum is also one of the few around the world that Art ' comes up with a result indicating the museum is permanently closed . This is in reference to the museum ' s old location , which was closed in 2013 . The new museum is , in fact , still open . Things to Do in Akita 6 6 Share on Facebook Share on X ( Twitter ) Copy link to share 1 Was this article helpful ? Help us improve the site Give Feedback By Victoria Kamila Community writer Following Follow +6 Information 1-4-2 Nakadori , Akita City , Akita , Japan ( Directions ) +81 18 853 8686 akita-museum-of-art.jp Book your trip Hotels Trip Find a nearby hotel Explore some of the best nearby hotels and deals for the perfect place to stay in Japan . Your stay dates : Search with Tripadvisor Start your trip now When do you want to travel ? My dates are flexible Get started Explore nearby Hotel Richmond Akita Eki-mae By Bonson Lam Kagura Restaurant at Hotel Richmond By Bonson Lam Akita Design Hub and Craft Center By Bonson Lam Kubota Castle By Hitoshi Kawai Discover Akita City Top Articles Recommended Popular 1 Freshness Burger Okinawa Fair Food 2 Guided Tour Foujita ' s , he became a patron of his . The free gallery located on the first floor displays temporary exhibitions of art and other projects or collaborations by local artists . Photography is not permitted in any of the galleries . Outside of the museum there is also a public market selling kitchenwares , teapots , those lucky waving cats , and those weird bear looking forest spirit things you see outside houses and temples all the time in Japan . It ' s also important to mention that a quick Google search for ' Akita Museum of known as Leonard Foujita—a name that was adapted after spending a considerable amount of his life in France and eventually acquiring French citizenship . He did exhibit a fondness for Akita in particular , artistically expressing this in a twenty meter long mural titled “ Annual Events in Akita ” . The mural depicts life in Akita throughout the four seasons and is given permanent status in this gallery . Also on display in separate galleries is the work of Hirano Masakichi , who was born in Akita . Having developed an affection for art , particularly in regards to claim about Google saying it is closed . Wikipedia , Google maps , and the official website all say it is open ( in English and in Japanese ) . Perhaps there a few outdated websites that list it as closed . I used to go to the old one when I was in Akita 2011-2012 . The new one with poolside cafe looks nice . Reply Show all 0 replies 1 comment in total houses a large collection of work by Fujita Tsugaru , arguably one of the most famous Western painters of Japan . The Tsugaru gallery is located up the free-standing staircase which is illuminated by a large triangular skylight and just past the cafe on the second floor . You do need to pay an entrance fee to view the gallery ( 300 yen ) , though the cafe and the museum shop are free to enter . The cafe itself features a great view accented by an infinity pool running along the glass-walled windows . Internationally , Fujita is also for Foreign Visitors to Experience Traditional Japanese Performing Arts and Learn About the History of Ginza ! News - Sponsored 3 Introducing Nomad eSIM Planning 4 Sakura-tei : The World ' s Largest Okonomiyaki Restaurant Tokyo 5 Kamakura : A Memorable Summer Getaway Kanagawa 6 A Foot Bath Path In Kaminoyama Yamagata 7 Fairytale Experiences at Nukumori-no-mori Shizuoka 8 A Comprehensive Guide to Nomad ' s Japan eSIM Planning 9 Tokyo Ginza Noh Week Tokyo - Sponsored 10 Fūrin , the Symbol of Japanese Summer Culture 1 A Guide to Japanese Visas Planning 2 Guide to Bringing Medicines Into Japan Planning

**GPT-4o Summarize: 210 yen** ✓

Figure 19: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

Round3 Summarization

### Qwen2-VL

**Full-page Screenshot:**



**Content:**

Reviews US \$ 66.00 View Top Restaurant Picks Near Akita Museum of Modern Art 1 . Bar Pasaporte Address : 204-1 Aza Takuboshita Fuke Otsutsumi Distance : 445m Bar Pasaporte No reviews yet Other Cuisine View 2 . Kuidoraku Price : \$ 8.00 Address : 7-2 Ekimaecho , Yokote , Akita 013-0036 Distance : 2.28km Kuidoraku No reviews yet Bars/Bistros US \$ 8.00 View 3 . Korakuen Yokoteten Price : \$ 5.00 Address : It is 28-1 , Sanmaibashi in Maego , Yokote-shi , Akita character Distance : 2.03km Korakuen Yokoteten No reviews yet Other Chinese Cuisine US \$ 5.00 View 4 . Ganso Kamiya Yakisoba Restaurant Address : Nakano-117-67 Oyashinmachi , Yokote , Akita 013-0051 , Japan Distance : 1.17km Ganso Kamiya Yakisoba Restaurant No reviews yet Fast Food View Verified Reviews of Akita Museum of Modern Art 第二号爱人 : 遇上秋田杆灯季 , 还是挺热闹的一个节日在美术馆里面的话 , 也有这种节日的气氛 , 氛围都还是相当不错的 , 而且的话里面虽然没有特别多的名画名品 , 但是因为 是免费进入 , 所以值得参观。Jedy Tan : 属于小众景点了 , 秋田县本来就不大 , 美术馆还在一个小 市里。但参观过能看出日本人近现代对西洋艺术的崇尚 xiaomoufa : 哈哈 , 这是一个非常美丽的美术馆 , 特别有意思的一个 景点。E30 \* \* \* 67 : 蛮大的美术馆 , 里面画有些我们国家古代的味道 , 蛮不错 的 Also Popular With Visitors to Akita Museum of Modern Art 1 . Sendai Umino-Mori Aquarium Price : \$ 14.30 Discount : \$ 2.04 Recommended sightseeing time : : 4-5 hours Address : Japan , 〒983-0013 and free . everything related to the ninjas were very fun . Edo Wonderland Nikko Edomura 4.5 / 5 43 Reviews No.3 of Best Things to Do in Nikko Theme Parks From US \$ 37.43 View Contents Akita Museum of Modern Art Opening Times Akita Museum of Modern Art Address Suggested Visit Duration for Akita Museum of Modern Art Featured Accommodation Near Akita Museum of Modern Art 1 . Hotel Plaza Annex Yokote 2 . Yokote Central Hotel 3 . Quad Inn Yokote 4 . Yokote Plaza Hotel Top Restaurant Picks Near Akita Museum of Modern Art 1 . Bar Pasaporte Hotels Flights Trains Cars Car Rentals Airport Transfers App Customer Support USD Search Bookings Sign in / Register Travel with Trip.com Travel Guide for Akita Museum of Modern Art in September ( Updated 2024 ) Akita Museum of Modern Art Opening Times Year round : 9:30-17:00 Akita Museum of Modern Art Address 62-46 , Akasaka Tomigazawa | Inside Akita Furusato Mura , Yokote , Akita Prefecture Suggested Visit Duration for Akita Museum of Modern Art 1-2 hours Featured Accommodation Near Akita Museum of Modern Art 1 . Hotel Plaza Annex Yokote Address 2 . Kuidoraku 3 . Korakuen Yokoteten 4 . Ganso Kamiya Yakisoba Restaurant Verified Reviews of Akita Museum of Modern Art Also Popular With Visitors to Akita Museum of Modern Art 1 . Sendai Umino-Mori Aquarium 2 . Tsugaru-han Neputa mura Village 3 . Suntopia World 4 . Edo Wonderland Nikko Edomura Contents Akita Museum of Modern Art Opening Times Akita Museum of Modern Art Address Suggested Visit Duration for Akita Museum of Modern Art Featured Accommodation Near Akita Museum of Modern Art 1 . Hotel Plaza Annex Yokote 2 . Yokote Central Hotel 3 . Quad Inn Yokote 4 . Yokote Plaza Hotel Top Restaurant Picks Near Akita Museum of Modern Art 1 . Bar Pasaporte 2 . Kuidoraku 3 . Korakuen Yokoteten 4 . Ganso Kamiya Yakisoba Restaurant Verified Reviews of Akita Museum of Modern Art Also Popular With Visitors to Akita Museum of Modern Art 1 . Sendai Umino-Mori Aquarium 2 . Tsugaru-han Neputa mura Village 3 . Suntopia World 4 . Edo Wonderland Nikko Edomura Popular Travelogues Bangkok Travelogue | Manila Travelogue | Tokyo Travelogue | Taipei Travelogue | Hong Kong Travelogue | Seoul Travelogue | Kuala Lumpur Travelogue | Los Angeles Travelogue | Shanghai Travelogue \$ 14.30 View 2 . Tsugaru-han Neputa mura Village Price : \$ 3.38 Address : Japan , 〒036-8332 Aomori , Hirosaki , Kamenokomachi , 6 1 Distance : 0.63 mi Tsugaru-han Neputa mura Village No reviews yet From US \$ 3.38 View 3 . Suntopia World Price : \$ 8.85 Recommended sightseeing time : : 0.5-1 day Address : 1-1 Kubo , Agano , Niigata 959-2212 , Japan Distance : 5.06 mi What travelers say : M515shunyi1618 : Very suitable for family outings of playground , Ferris wheel , pirate ship are fun . Suntopia World 5 / 5 9 Reviews Amusement Ekimaecho Price : \$ 53.00 Distance : 2.23km Hotel Plaza Annex Yokote 4.3 / 5 37 Reviews -5 % US \$ 53.00 View 2 . Yokote Central Hotel Address : Heiwacho 9-10 Price : \$ 50.00 Distance : 2.91km Yokote Central Hotel 3.7 / 5 18 Reviews -5 % US \$ 50.00 View 3 . Quad Inn Yokote Address : Sekibata-52-1 Yasuda Price : \$ 58.00 Distance : 1.9km Quad Inn Yokote 4.1 / 5 14 Reviews US \$ 58.00 View 4 . Yokote Plaza Hotel Address : 7-1 Ekimaecho Price : \$ 66.00 Distance : 2.27km Yokote Plaza Hotel 3.9 / 5 10

Qwen2-VL Summarize: US \$ 66.00

✘

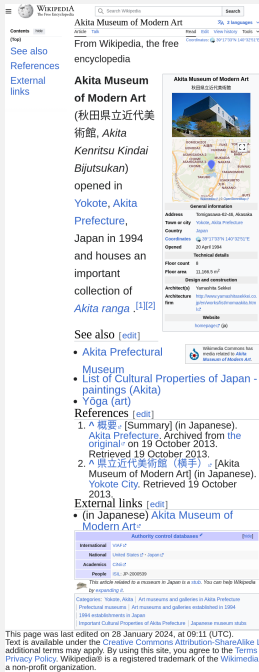
Figure 20: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

**Round3 Summarization**

### LLaVA-OneVision

**Full-page Screenshot:**



**Content:**

Akita Museum of Modern Art 秋田県立近代美術館 横手 秋田県立近代美術館 (横手) [ Akita Museum of Modern Art ] ( in Japanese ) . Yokote City . Retrieved 19 October 2013 . External links [ edit ] ( in Japanese ) Akita Museum of Modern Art hide Authority control databases International VIAF National United States Japan Academics CiNii People ISIL : JP-2000539 This article related to a museum in Japan is a stub . You can help Wikipedia by expanding it . v t e Retrieved from " https : //en.wikipedia.org/w/inde x.php ? title=Akita\_Museum\_of\_Modern\_Art & oldid=1199934740 " Categories : Yokote , Akita Art museums and galleries in Akita Prefecture Prefectural museums Art museums and galleries established in 1994 1994 establishments in Japan Important Cultural Properties of Akita Prefecture Japanese museum stubs Hidden categories : Pages using gadget WikiMiniAtlas CS1 uses Japanese-language script ( ja ) CS1 Japanese-language sources ( ja ) Use dmy dates from November 2019 Articles with short description Short description is different from Wikidata Infobox mapframe without OSM relation ID on Wikidata Coordinates on Wikidata Articles containing Japanese-language text Commons category link is on Wikidata Articles with Japanese-language sources ( ja ) All stub articles Pages using the Kartographer extension

**LLaVA-OneVision Summarize: free** ✘

Figure 21: Response and middle results comparison of GPT-4o (OpenAI, 2024b), Qwen2-VL-7B (Qwen Team, 2024), and LLaVA-OneVision-7B (Li et al., 2024b) in the end-to-end task.