

# M<sup>3</sup>Searcher: Modular Multimodal Information Seeking Agency with Retrieval-Oriented Reasoning

Anonymous ACL submission

## Abstract

Recent advances in DeepResearch-style agents have demonstrated strong capabilities in autonomous information acquisition and synthesis from real-world web environments. However, existing approaches remain fundamentally limited to text modality. Extending autonomous information-seeking agents to multimodal settings introduces critical challenges: the specialization-generalization trade-off that emerges when training models for multimodal tool-use at scale, and the severe scarcity of training data capturing complex, multi-step multimodal search trajectories. To address these challenges, we propose M<sup>3</sup>Searcher, a modular multimodal information-seeking agent that explicitly decouples information acquisition from answer derivation. M<sup>3</sup>Searcher is optimized with a retrieval-oriented multi-objective reward that jointly encourages factual accuracy, reasoning soundness, and retrieval fidelity. In addition, we develop MM-SearchVQA, a multimodal multi-hop dataset to support retrieval centric RL training. Experimental results demonstrate that M<sup>3</sup>Searcher outperforms existing approaches, exhibits strong transfer adaptability and effective reasoning in complex multimodal tasks.

## 1 Introduction

DeepResearch-style agents have recently demonstrated striking proficiency in acquiring and synthesizing information from real-world web environments, as exemplified by OpenAI DeepResearch (OpenAI, 2025) and Gemini DeepResearch (Google, 2025). These advances have spurred a growing research effort to equip large language models (LLMs) with reasoning-intensive search capabilities (Shao et al., 2025). Most approaches leverages reinforcement learning (RL) to train models to interact with web search engines (e.g. Google

Search), planning, gathering and synthesizing information through multi-step deliberation (Jin et al., 2025; Zheng et al., 2025). However, these approaches remain confined to text modality, even though real-world user information needs are inherently multimodal (e.g. visual perception).

Extending autonomous information-seeking agents to multimodal inputs is therefore an essential step for building general intelligent systems. Nevertheless, this transition introduces several fundamental challenges: (i) **Specialization-Generalization Trade-off:** Training models to internalize multimodal tool-use policies comes at the expense of general reasoning capacity (Kalajdzievski, 2024; Li et al., 2024a), yet the effectiveness of multimodal RAG systems critically relies on a backbone model whose core reasoning performance remains robust and uncompromised. (ii) **Training Data Scarcity:** Existing datasets that capture complex, multi-step search trajectories are primarily designed for evaluation purposes (Wei et al., 2025), whereas large-scale training corpora provide only shallow reasoning paths, such as InfoSeek (Chen et al., 2023). This discrepancy hinders models from developing long-horizon information-seeking strategies.

To resolve these challenge, and inspired by the modular design of Jiang et al. (2025), we decouple the information-seeking process from answer derivation. Specifically, we introduce a lightweight and trainable MLLM, termed **M<sup>3</sup>Searcher**, that serves as a dedicated modular multimodal information seeking agency. Its role is to execute a mulitmodal reasoning-intensive information seeking process (Shao et al., 2025). Specifically, it interprets non-textual inputs (e.g. visual recognition, OCR) and dynamically coordinating search strategies across heterogeneous modalities to assemble comprehensive and contextually relevant evidence. The gathered information is subsequently provided to a downstream answer generator, which performs reasoning over the curated evidence and formulates

<sup>2</sup>✉ Corresponding author.

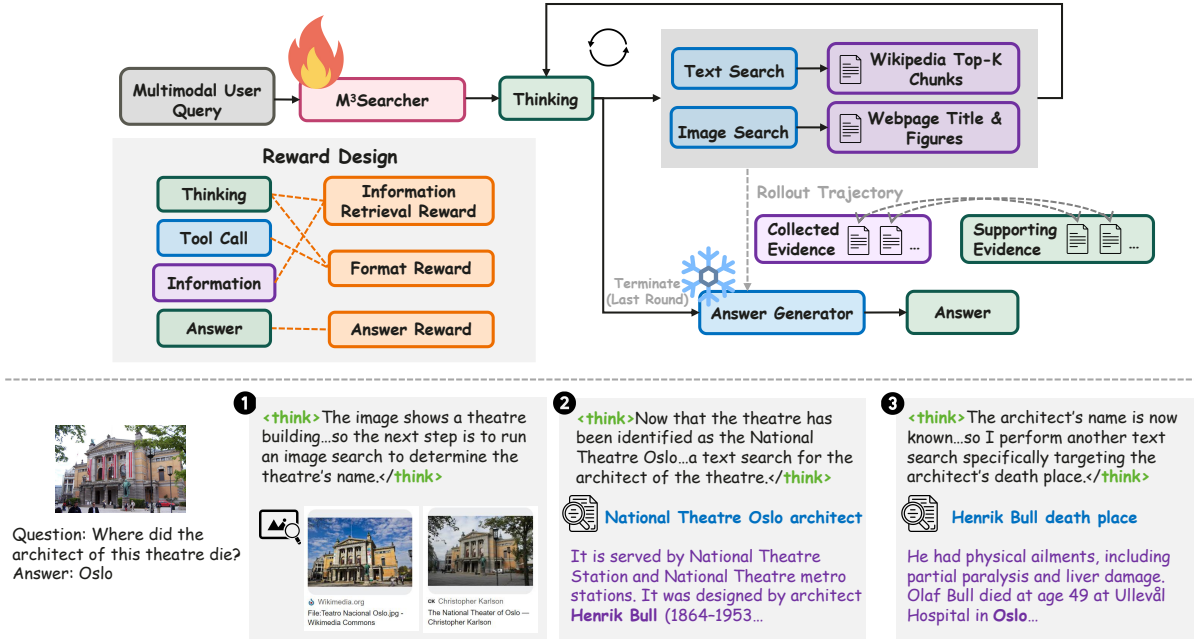


Figure 1: The architecture of M<sup>3</sup>Searcher.

the final response to the user query. To effectively train M<sup>3</sup>Searcher, we further propose a decoupled reinforcement learning framework with the following contributions:

- 1. Dataset Construction:** We introduce **MM-SearchVQA**, a dataset demanding rigorous multimodal information seeking. Each instance enforces answer uniqueness and is accompanied by automatically extracted supporting evidence. By encompassing a broad spectrum of domains, difficulty levels and search intensities, the dataset encourages the model to learn the distinct control policies required for determining *when* to search, *what* to query, and *how* to integrate external knowledge.
- 2. Decoupled Multimodal Information Seeking:** M<sup>3</sup>Searcher focuses exclusively on optimizing heterogeneous search scheduling for maximizing information acquisition. To realize this, we introduce a specialized "expert answer generator" tool, which is triggered only once the context is deemed sufficient and well-grounded for the following reasoning process. This modularity allows the search strategies to remain highly adaptive while maintaining the reasoning capacity of a robust backbone within the MRAG system. It also renders the generator modality-agnostic, accommodating both pure textual LLMs (e.g., DeepSeek-R1) and MLLMs (e.g., GPT-4o).

- 3. Retrieval-Oriented Multi-Objective Reward:** We employ a multi-objective reward modeling framework that jointly optimizes answer accuracy, reasoning validity, and retrieval quality. To ensure that the model genuinely grounds its inferences in retrieved evidence rather than exploiting spurious shortcuts, we incorporate a retrieval reward that evaluates the completeness of textual information gathering and the accuracy, and interpretive soundness of visual reasoning.

We conduct comprehensive evaluations of M<sup>3</sup>Searcher across real-world benchmarks to assess its effectiveness. M<sup>3</sup>Searcher outperforms both prompt-engineered agents and end-to-end trained counterparts. Moreover, it exhibits strong robustness and adaptability, as evidenced by stable performance under multiple transfer scenarios involving variations in search engines and answer generators.

## 2 M<sup>3</sup>Searcher

### 2.1 Task Formulation

We consider a multimodal query  $(v, q)$  where  $v$  is the visual component and  $q$  is the textual component, with its ground-truth answer  $a$ . A trainable MLLM is formalized as an information-seeking agent  $\mathcal{F}$ , which engages in iterative interaction with a multimodal tool set  $\mathcal{T}$ . Through a sequence of tool invocations and intermediate reasoning steps,

139 the agent incrementally acquires task-relevant ev- 186  
140 idence and integrates the retrieved information to 187  
141 derive the final answer: 188

$$142 \mathcal{F}(q, v, \mathcal{T}) \rightarrow a. \quad (1) \quad 189$$

## 143 2.2 Decoupled Agentic MRAG 192

144 Existing MRAG approaches either leverage a large- 193  
145 scale MLLM with elaborate prompt engineering, 194  
146 or train a smaller MLLM end-to-end (Geng et al., 195  
147 2025; Wu et al., 2025b). This dichotomy introduces 196  
148 a fundamental dilemma. Large-scale MLLMs (e.g. 197  
149 GPT-4o) excel in emergent reasoning, but lack op- 198  
150 timization for real-world web integration. Con- 199  
151 versely, smaller models (e.g., Qwen2.5-VL-7B) 200  
152 optimized for web search exhibit a concurrent 201  
153 degradation in general reasoning, limiting their 202  
154 utility as a primary backbone for the overall sys- 203  
155 tem. To address this dilemma, and drawing in- 204  
156 spiration from modular architectures (Jiang et al., 205  
157 2025), we introduce a decoupled MRAG architec- 206  
158 ture that separates the information-seeking process 207  
159 from the answer generation. As depicted in Figure 2, M<sup>3</sup>Searcher is solely responsible for compre- 208  
160 hending multimodal queries, formulating iterative 209  
161 search strategies across heterogeneous search tools, 210  
162 and determining the optimal search termination 211  
163 point. The final evidential data is then passed to 212  
164 a dedicated answer generation for synthesis. This 213  
165 architectural decoupling offers two key advantages: 214  
166 it preserves the reasoning fidelity of the large-scale 215  
167 backbone while allowing for targeted optimization 216  
168 of information-seeking capabilities, and simultane- 217  
169 ously removes modality constraints, enabling the 218  
170 use of modality-agnostic generators (e.g., GPT-4o, 219  
171 DeepSeek-R1). 220  
172

## 173 2.3 Multimodal Tools Implementation 225

174 We equip M<sup>3</sup>Searcher with three essential tools: 226  
175 an image search tool, a text search tool, and an 227  
176 answer generator tool. To enable effective RL, we 228  
177 developed a stable and high-concurrency tool en- 229  
178 vironment. For image search, we integrate the 230  
179 Serper API<sup>1</sup> to perform reverse image retrieval. 231  
180 Given an input image, the API returns visually 232  
181 similar images, together with their corresponding 233  
182 website titles and URLs. Since the Serper API 234  
183 produces highly stable results, we incorporate a 235  
184 caching mechanism to reduce resource consump-  
185 tion and accelerate the search process. For text

<sup>1</sup><https://serpapi.com/>

186 search, we utilize the 2025 wikipedia dump<sup>2</sup> as 187  
188 knowledge source. A retrieval–reranking pipeline, 189  
190 built upon the E5 models (Wang et al., 2022), is 191  
192 used to retrieve semantically relevant document 193  
194 chunks given a user query. Finally, the answer 195  
196 expert tool employs a high-capacity LRM which 197  
198 consumes the trajectory of information-seeking and 199  
200 synthesizes a final response. 201

## 202 2.4 Decoupled Multi-turn Rollout 203

204 M<sup>3</sup>Searcher processes the query through three 205  
206 core operational states. In the **Think** state, the 207  
208 model conducts a fine-grained inspection of the 209  
210 visual component  $v$  and performs contextual in- 211  
212 ference across modalities, integrating visual and 213  
214 textual cues to construct a coherent situational un- 215  
216 derstanding. When additional information is re- 217  
218 quired, M<sup>3</sup>Searcher transitions to the **Tool\_Call** 219  
220 state, dynamically invoking external tools to re- 221  
222 trieve supplementary evidence from real-world 223  
224 sources. The retrieved outputs are then encap- 225  
226 sulated as **Information**, which re-enters the rea- 227  
228 soning loop to refine and expand the model’s un- 229  
230 derstanding. M<sup>3</sup>Searcher operates iteratively upon 231  
232 these three states, allowing for progressive refine- 233  
234 ment of its understanding and retrieval strategy. 234  
235 Formally, at each time step  $t$ , the M<sup>3</sup>Searcher ex- 236  
237 ecution can be represented as a tuple  $(\alpha_t, C_t, I_t)$ , 238  
239 where  $\alpha_t$  represents the reasoning process,  $C_t$  is 239  
240 the tool invocation, and  $I_t$  is the tool response. The 240  
241 full rollout trajectory can thus be expressed as: 241

$$242 \mathcal{T} = \{O_1, \alpha_1, C_1, I_1, \dots, O_t, \alpha_t, C_t, I_t\}. \quad (2) \quad 243$$

244 Under the decoupled agentic design, the final tool 245  
246 invocation of M<sup>3</sup>Searcher is required to invoke the 246  
247 answer expert. Consequently, the terminal tool 247  
248 response  $I_t$  provides the final answer to the user 248  
249 query  $q$ : 249  
250

$$251 I_t = \mathcal{F}(q). \quad (3) \quad 252$$

253 The prompt governing the rollout procedure is de- 254  
255 tailed in Appendix B. 255

## 256 2.5 Multi-Objective Rewrad Modeling 257

258 The goal of M<sup>3</sup>Searcher is to perform a multimodal, 259  
260 reasoning-intensive information-seeking process. 260  
261 It must progressively and comprehensively gather 261  
262 relevant evidence across multiple hops to support 262  
263 the downstream generator in generating accurate 263  
264

<sup>2</sup><https://dumps.wikimedia.org/>

answers. To achieve this, we formulate a multi-objective, retrieval-oriented RL reward function that jointly optimize accuracy, completeness and relevance of the information acquisition process.

**Format Reward** The format reward  $R_{format}$  enforces strict compliance with the syntactic and structural constraints specified in the prompt. For example, tool invocations are required to follow the correctly structured parsing format with valid parameterization; and the trajectory must terminate with a call to the answer generator tool. Any deviation from these requirements incurs a strong penalty of an absolute reward of -1.

**Answer Reward** The answer reward,  $R_{answer}$ , measures the semantic correctness of the final output  $I_t$  with respect to the reference solution. Rather than relying on brittle exact string matching, we employ an LLM-as-Judge evaluation strategy, which confers both flexibility and robustness in cases where multiple equivalent phrasings or semantically consistent answers are acceptable. The complete scoring prompt used for this evaluation is provided in the Appendix.

**Information Retrieval Reward** The information retrieval reward is designed to assess the fidelity and completeness of the information acquired to solve a multi-hop user query, independent of the capabilities of the downstream answer generator. The evaluation of this information acquisition is divided by modality.

For the visual modality, when processing retrieved images or relying on internal knowledge within the model’s training cutoff, it may exhibit three distinct behaviors during the **Think** state: (1) correctly identifying the key visual elements, (2) demonstrating uncertainty and refraining from explicit recognition while offering a descriptive interpretation, or (3) producing an incorrect recognition. To shape this behavior, we assign graded rewards  $R_{ImgRetrieval}$  of 0.5, 0.25, and 0, respectively. This reward structure encourages the model to adopt a more cautious and self-aware strategy when reasoning over visual inputs. For the textual modality, we assess the degree to which the information conveyed to the answer generator aligns with the reference evidence in MMSearchVQA. This metric quantifies whether M<sup>3</sup>Searcher identifies all necessary pieces of information required for solving the query, thereby enhancing information completeness and mitigating reasoning shortcuts

that may yield correct answers without genuine verification (e.g., a builder’s place of death is not always the same as the building’s location). The textual retrieval reward, denoted as  $R_{TextRetrieval}$ , is defined as a percentage score ranging from 0 to 0.5, representing the proportion of reasoning hops successfully supported by retrieved evidence. Specifically, we compare each reference evidence against both the **Information** and **Think** states. To ensure robust evaluation, we employ a LLM-as-Judger method to assess both modality reward score:

$$R_{ImgRetrieval} = LLM(\alpha_i), \quad (4)$$

$$R_{TextRetrieval} = LLM(\alpha_i, C_i). \quad (5)$$

The detailed judging prompt is provided in the Appendix. The final reward is:

$$R = R_{format} + R_{answer} + R_{Retrieve}, \quad (6)$$

$$R_{Retrieve} = R_{TextRetrieval} + R_{ImgRetrieval}. \quad (7)$$

## 2.6 RL Training

To enhance the model’s capability for information seeking and web-environment interaction within the MRAG framework, we adopt Group-Relative Policy Optimization (GRPO) (Shao et al., 2024). For each input multimodal question  $q$ , the current policy  $\pi_\theta$  samples a group of trajectories  $y_1, \dots, y_G$ . Then the optimization objective of GRPO is formulated as:

$$\mathcal{J}(\theta) = \mathbb{E}_{i,t} \min \left[ \rho_t^i A_t^i, \text{clip}(\rho_t^i, 1 - \epsilon, 1 + \epsilon) A_t^i \right] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta || \pi_{\text{ref}}], \quad (8)$$

where  $\rho_t^i$  represents the importance sampling ratio between the updated and previous policies and  $A_t^i$  is an estimator of the advantage at time step  $t$ :

$$A_t^i = \frac{R_i - \text{mean}(\{R_i\})}{\text{std}(\{R_i\})}. \quad (9)$$

The hyperparameter  $\beta$  controls the KL divergence penalty, constraining the deviation from the reference policy to ensure stable updates. The context for policies includes both model-generated outputs and tool responses. To prevent external knowledge sources from biasing policy learning, we apply a loss mask over all tool-response tokens. This ensures that policy gradients are computed exclusively for LLM-generated tokens, enabling precise

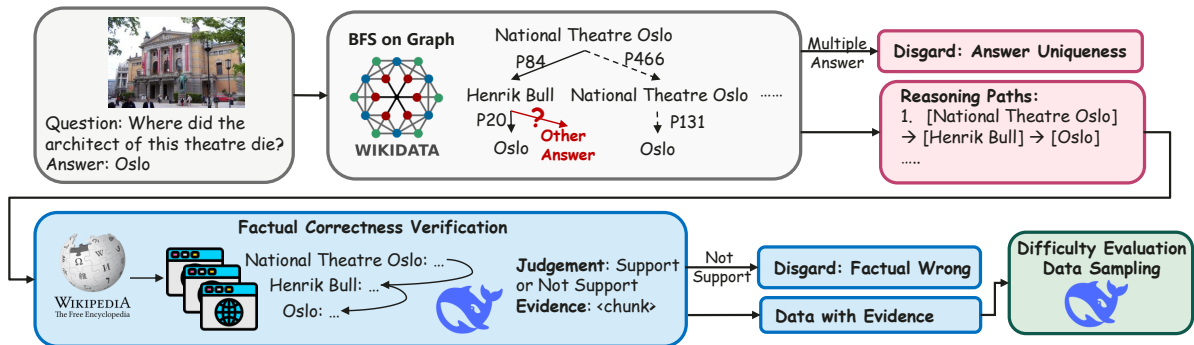


Figure 2: The MMSearchVQA data construction pipeline.

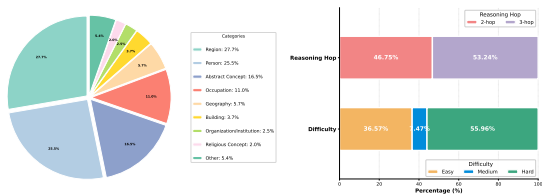


Figure 3: Overview of MMSearchVQA dataset statistics. The left figure summarizes the domain distribution, and the right figure reports the distribution of question difficulty levels and reasoning hops.

optimization of search planning and multimodal information-seeking capabilities within the MRAG system.

### 3 MMSearchVQA Dataset

Existing Visual Question Answering (VQA) datasets typically fall into two categories. Automatically constructed datasets (Chen et al., 2023; Cheng et al., 2025; Wu et al., 2025b; Fu et al., 2025), such as InfoSeek, involve shallow reasoning chains — often limited to two-hop queries solvable through a simple sequence of image search followed by text search. In contrast, manually curated datasets such as MM-BrowseComp (Li et al., 2025b) feature more complex, multi-step reasoning but are expensive and difficult to scale. To address this limitation, we introduce MMSearchVQA, a dataset designed to foster the development of models for advanced information-seeking reasoning. MMSearchVQA not only requires deeper search and reasoning but also provides explicit supporting evidence that underpin the reasoning and answering processes.

As illustrated in Figure 2, our dataset is constructed upon ReasonVQA (Tran et al., 2025), which is derived from the Wikidata. We first perform a BFS traversal on the Wikidata graph, iden-

tifying all potential reasoning chains associated with each question. During this traversal, we discard questions that yield multiple valid answers to ensure answer uniqueness, and we retain only those samples that require at least two reasoning hops. Following the extraction of candidate reasoning paths, we conduct cross-validation against Wikipedia using the DeepSeek models. Each reasoning hop, including the final answering, must be consistently supported by evidence drawn from relevant Wikipedia content to be both factually accurate and temporally valid. Instances that fail to meet these criteria are excluded. During this verification process, we also extract fine-grained supporting evidence from the corresponding Wikipedia passages for each reasoning step, thus enhancing the interpretability and traceability of the reasoning process. To further characterize the cognitive difficulty of the resulting dataset, we employ DeepSeek-V3 (Guo et al., 2025) to answer each question three times and categorize questions into three levels: easy (all correct), medium (partially correct), and hard (all incorrect). This procedure yields a principled estimate of reasoning complexity across samples. To cultivate an information seeker capable of performing deep and precise searches, we prioritize training data that exhibit deeper information needs and greater reasoning complexity. Accordingly, we downsample easy questions to half the number of hard examples, ensuring a balanced yet challenging dataset. In total, the curated dataset contains 6,000 questions, with comprehensive statistics presented in Figure 3.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets** We adopt MMSearchVQA dataset as the training corpus. We evaluate performance on

Table 1: The overall performance of M<sup>3</sup>Searcher compared with baseline approaches. MMSearchVQA is in-domain benchmark with wikipedia based search and other benchmarks are transferred to Google Search Engine. The best performance is highlighted in **bold**, and the second-best performance is underlined.

Method	Backbone	In-Domian (Wiki Search)	→ Out-Domain (Google Search)		
		MMSearchVQA	InfoSeek	MMSearch	MRAG-Bench
<b>No Agency</b>					
Direct	Qwen3-VL-235B-A22B	40.42	40.16	28.23	9.22
	Qwen2.5-VL-72B	43.12	35.22	15.29	14.12
	Qwen2.5-VL-7B	31.12	23.50	11.69	8.57
RAG	Qwen3-VL-235B-A22B	29.79	31.75	30.83	<b>33.67</b>
	Qwen2.5-VL-72B	40.37	33.26	<u>46.15</u>	20.08
	Qwen2.5-VL-7B	30.00	31.75	38.09	15.68
<b>Prompt Engineered Agents</b>					
OmniSearch	Qwen2.5-VL-72B	45.65	40.60	15.00	27.07
	Qwen2.5-VL-7B	22.91	25.17	22.22	23.96
CogPlanner	Qwen2.5-VL-72B	48.37	41.72	39.77	29.12
	Qwen2.5-VL-7B	22.12	26.22	27.48	28.23
<b>End-to-End Agents</b>					
MMSearch-R1-7B <sub>Retrain</sub>		31.63	20.20	7.02	27.60
MMSearch-R1-7B <sub>Release</sub>		20.50	37.06	12.28	19.20
<b>Decoupled Agents w/o Training</b>					
Qwen3-30B-A3B		31.50	31.20	36.69	24.20
Qwen2.5-VL-7B		34.12	33.80	36.09	27.20
<b>M<sup>3</sup>Searcher</b>					
Qwen3-30B-A3B		54.75	39.61	55.62	24.91
→ Transfer LRM Answer Generator					
DeepSeek-V3		56.87	40.33	60.95	29.12
DeepSeek-R1		59.25	<b>42.50</b>	<b>63.30</b>	<u>30.00</u>
→ Transfer MLLM Answer Generator					
Qwen2.5-VL-7B		57.00	39.44	61.54	19.95
Qwen2.5-VL-72B		<b>59.50</b>	40.20	59.17	27.20

both in-domain and out-of-domain benchmarks. The in-domain evaluation uses the MMSearchVQA test set, while out-of-domain evaluation is conducted on three publicly available VQA datasets: MMSearch (Jiang et al., 2024), Infoseek (Chen et al., 2023), and MRAG-Bench (Hu et al., 2024).

**Baselines and Metrics** We compare M<sup>3</sup>Searcher against four categories of methodologies: (1) No-agency: We directly prompt MLLMs and use a fixed RAG pipeline comprising image retrieval, query rewriting, text retrieval, and answer generation. (2) Prompt-engineered agents: We select OmniSearch (Li et al., 2024b) and CogPlanner (Yu et al., 2025), both of which coordinate multiple agents via hand-crafted prompts for multimodal reasoning and retrieval. (3) End-to-end agents: We include MMSearch-R1 (Wu et al., 2025b) as a representative method optimized through end-to-end RL training. (4) Decoupled agents without specialized training: We employ a decoupled architecture

in which Qwen2.5-VL-7B is used for information-seeking operations, without any tuning. We adopt LLM-as-Judge as the evaluation metric, which is well-aligned with the answer accuracy reward.

**Transfer Experiment Settings** For M<sup>3</sup>Searcher, we conduct two sets of transfer experiments: (1) Search engine transfer: For the MMSearchVQA benchmark, which is built on Wikipedia-based content, we employ our in-house text search tool (described in Section 2.3) to retrieve the top 10 most relevant text chunks. For other benchmarks based on open-domain web data, we switch to Google Search via the Serper API<sup>3</sup>, also keeping the top 10 retrieved results. (2) Answer generator transfer: We explore the transfer of answer generators by incorporating models from both the DeepSeek series (Guo et al., 2025; Liu et al., 2024) and the Qwen-VL series (Bai et al., 2025).

<sup>3</sup><https://serper.dev/>

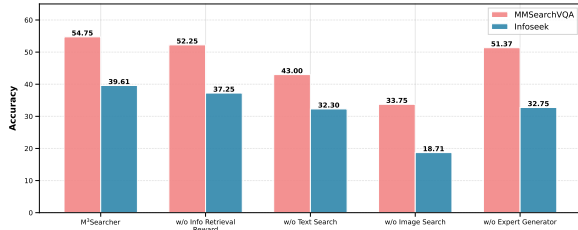


Figure 4: Ablation study.

**Implementation Details** For the baseline methodologies, we adopt Qwen3-VL-30B-A3B, Qwen2.5-VL-72B, and Qwen2.5-VL-7B as backbone models. Specifically, MMSearch-R1 employs Qwen2.5-VL-7B for both end-to-end training and inference. For M³Searcher, we utilize Qwen2.5-VL-7B as the trainable planner and Qwen3-30B-A3B (Yang et al., 2025) as the answer generator during training. Verl (Sheng et al., 2024) is used for multi-turn RL training.

## 5 Main Results

The overall performance of M³Searcher across multiple benchmarks is summarized in Table 1. Several insights can be drawn from these results: (1) **Baseline methodologies exhibit unstable variability in performance across benchmarks.** In most cases, either fixed RAG pipelines or prompt-engineered agents attain the strongest results. This suggests that explicit, hand-crafted prompt engineering provides a competitive advantage that decoupled, untrained agents fail to surpass. Agents with specialized training display inconsistent performance and unstable generalization: for example, MMSearch-R1 performs competitively on Infoseek, but its performance drops sharply on out-of-distribution tasks. (2) **M³Searcher demonstrates robust and strong performance across various generalization and transfer settings.** It provides high-quality, correctly excavated evidence and both multimodal and purely textual backbones can reliably synthesize accurate answers, lifting the modality constraints. Notably, DeepSeek-R1 answer generator emerges as the top performer, underscoring the critical role of the inherent reasoning capability of the backbone model in the overall MRAG system effectiveness. M³Searcher also maintains stable performance under search-engine transfer, exhibiting no degradation when switching search tools. This robustness highlights its high robustness to variations in the underlying information source, and further indicates that a self-built textual search

engine is fully sufficient for on-policy RL training — particularly important given the prohibitive cost of commercial search engines.

## 6 Analysis

**Ablation study.** We evaluate the contribution of each core component in M³Searcher by removing the Information Retrieval Reward, the image search tool, the text search tool and the answer generator tool. As shown in Figure 4, each component provides a measurable performance gain, underscoring their collective importance to the overall system effectiveness.

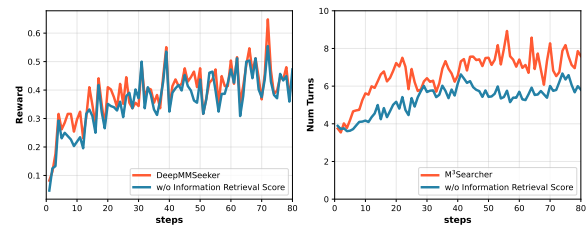


Figure 5: Training dynamics of reward and rollout turn counts with and without the information-retrieval reward.

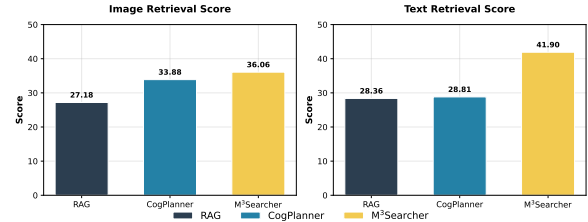


Figure 6: Text retrieval score and image retrieval score of M³Searcher compared with CogPlanner and RAG pipeline baselines.

**Retrieval-oriented rewards enhance the breadth and completeness of information seeking.** To rigorously evaluate the impact of retrieval-oriented reward design, we analyze the training dynamics presented in Figure 5. The results indicate that incorporating an information retrieval reward leads to consistently higher reward signals. Consequently, M³Searcher engages in a greater number of information-seeking turns. This enables a broader coverage of relevant information and yielding final evidence that is both more complete and reliable, as demonstrated in Figure 6.

**RL enhances heterogeneous tool coordination and improves the model’s ability to leverage the**

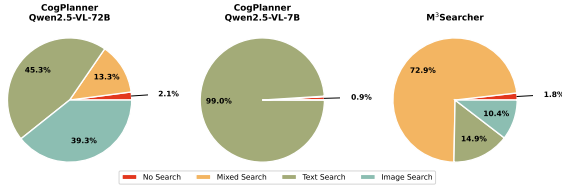


Figure 7: Tool usage statistics on MMSearchVQA.

**image-search tool** We analyze tool usage patterns on the MMSearchVQA benchmark in Figure 7. The results reveal a pronounced bias in the pre-trained Qwen2.5-VL-7B model as it overwhelmingly favors text search tool while almost never invoking the image search tool. After RL this imbalance is substantially mitigated. It invokes a more diverse mixture of search actions, with a notably increased reliance on the image-search tool, indicating that RL helps the model internalize when visual external information is necessary for successful reasoning.

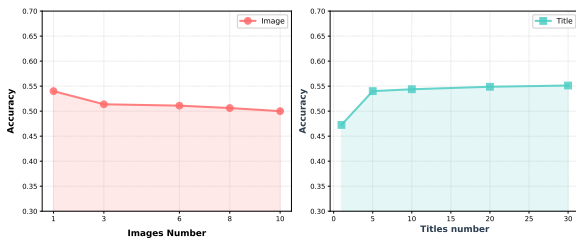


Figure 8: Performance with varying numbers of returned relevant images and associated webpage titles in the image search tool.

**The Design of Image Search Tools as a Core Factor in Agentic MRAG Performance** Our empirical observations indicate that the design of image search tools constitutes a critical determinant of performance. Specifically, we analyzed the influence of the returned relevant images and their associated webpage titles. As depicted in Figure 8, increasing the number of returned images lead to performance degradation, potentially attributed to the redundancy or near-similarity of multiple image inputs, which may introduce noise or confusion into the model’s visual feature extraction process. Conversely, increasing the volume of textual returned information (webpage titles) demonstrates a positive correlation with performance since the webpage titles provide crucial context for interpreting the visual query. Based on this quantitative analysis, we adopt a design choice for the image search tool that returns the top-1 image along with

the top-30 associated webpage titles, which balances informative context with minimal visual redundancy.

## 7 Related Work

With the introduction of DeepResearch by leading AI organizations, including OpenAI (OpenAI, 2025), Google (Google, 2025), and Perplexity (Perplexity, 2025), these systems have demonstrated strong potential in solving complex multi-step reasoning tasks. Recent advances highlight reinforcement learning (RL) as a promising paradigm and OpenAI’s technical report explicitly demonstrates the effectiveness of employing RL to strengthen the multi-step decision-making and retrieval abilities (Jaech et al., 2024). Notable works such as SearchR1 (Jin et al., 2025; Song et al., 2025) mark an early milestone by incorporating web-search tool interaction into textual question-answering scenarios, achieving substantial performance gains. Following this, the Web Agents series developed by the Qwen team (Li et al., 2025a; Wu et al., 2025a) further optimizes information-seeking behaviors in complex, non-linear reasoning tasks. However, despite these advances, rare attention has been given to the optimization of MRAG systems, where reasoning must integrate and synthesize heterogeneous modalities. Existing work (Geng et al., 2025; Wu et al., 2025b; Narayan et al., 2025) employ an end-to-end RL paradigm for VQA tasks, which inadvertently restrict the MRAG backbone to relatively small models (e.g., Qwen2.5-VL-7B). This constraint imposes a substantial performance ceiling, limiting the practical effectiveness of these systems in real-world deployments.

## 8 Conclusion

We present M³Searcher, a lightweight and trainable multimodal information seeker that decouples retrieval from answer generation in MRAG systems. By focusing on adaptive, reasoning-intensive search over heterogeneous sources, M³Searcher preserves the reasoning capacity of downstream generators while efficiently aggregating contextually relevant evidence. Experiments on MMSearchVQA and real-world benchmarks demonstrate strong performance and robustness across different search engines and generators.

## 567 Limitations

568 We discuss several limitations of M<sup>3</sup>Searcher as  
569 follows. First, although M<sup>3</sup>Searcher adopts a mod-  
570 ular architecture, its effectiveness is inherently con-  
571 strained by the scale and diversity of the available  
572 tool set. Extending the agent to operate over a  
573 broader and more heterogeneous collection of real-  
574 world tools would substantially enlarge the action  
575 space and increase planning complexity. Second,  
576 while MMSearchVQA facilitates retrieval-centric  
577 multimodal training, the constructed queries are  
578 predominantly characterized by relatively long rea-  
579 soning trajectories compared to those in existing  
580 training corpora. More complex scenarios that re-  
581 quire substantially deeper multi-step search and  
582 decision-making processes remain underexplored.  
583 Extending the dataset construction pipeline there-  
584 fore represents an important direction for future  
585 research.

## 586 References

587 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
588 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie  
589 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl  
590 technical report. [arXiv preprint arXiv:2502.13923](#).

591 Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, So-  
592 ravit Changpinyo, Alan Ritter, and Ming-Wei Chang.  
593 2023. Can pre-trained vision and language models  
594 answer visual information-seeking questions? [arXiv](#)  
595 [preprint arXiv:2302.11713](#).

596 Xianfu Cheng, Wei Zhang, Shiwei Zhang, Jian Yang,  
597 Xiangyuan Guan, Xianjie Wu, Xiang Li, Ge Zhang,  
598 Jiaheng Liu, Yuying Mai, and 1 others. 2025. Sim-  
599 plevqa: Multimodal factuality evaluation for mul-  
600 timodal large language models. [arXiv preprint](#)  
601 [arXiv:2502.13059](#).

602 Mingyang Fu, Yuyang Peng, Benlin Liu, Yao Wan, and  
603 Dongping Chen. 2025. Livevqa: Live visual knowl-  
604 edge seeking. [arXiv preprint arXiv:2504.05288](#).

605 Xinyu Geng, Peng Xia, Zhen Zhang, Xinyu Wang,  
606 Qiuchen Wang, Ruixue Ding, Chenxi Wang, Jia-  
607 long Wu, Yida Zhao, Kuan Li, and 1 others.  
608 2025. Webwatcher: Breaking new frontiers of  
609 vision-language deep research agent. [arXiv preprint](#)  
610 [arXiv:2508.05748](#).

611 Google. 2025. Gemini deep research system card.

612 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao  
613 Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-  
614 rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.  
615 Deepseek-r1: Incentivizing reasoning capability in  
616 llms via reinforcement learning. [arXiv preprint](#)  
617 [arXiv:2501.12948](#).

Wenbo Hu, Jia-Chen Gu, Zi-Yi Dou, Mohsen Fayyaz,  
Pan Lu, Kai-Wei Chang, and Nanyun Peng. 2024.  
Mrag-bench: Vision-centric evaluation for retrieval-  
augmented multimodal models. [arXiv preprint](#)  
[arXiv:2410.08182](#). 618 619 620 621 622

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-  
son, Ahmed El-Kishky, Aiden Low, Alec Helyar,  
Aleksander Madry, Alex Beutel, Alex Carney, and 1  
others. 2024. Openai o1 system card. [arXiv preprint](#)  
[arXiv:2412.16720](#). 623 624 625 626 627

Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanmin Wu,  
Jiayi Lei, Pengshuo Qiu, Pan Lu, Zehui Chen, Guan-  
glu Song, Peng Gao, and 1 others. 2024. Mm-  
search: Benchmarking the potential of large mod-  
els as multi-modal search engines. [arXiv preprint](#)  
[arXiv:2409.12959](#). 628 629 630 631 632 633

Pengcheng Jiang, Xueqiang Xu, Jiacheng Lin, Jin-  
feng Xiao, Zifeng Wang, Jimeng Sun, and Jiawei  
Han. 2025. s3: You don't need that much data  
to train a search agent via rl. [arXiv preprint](#)  
[arXiv:2505.14146](#). 634 635 636 637 638

Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon,  
Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei  
Han. 2025. Search-r1: Training llms to reason and  
leverage search engines with reinforcement learning.  
[arXiv preprint arXiv:2503.09516](#). 639 640 641 642 643

Damjan Kalajdzievski. 2024. Scaling laws for forget-  
ting when fine-tuning large language models. [arXiv](#)  
[preprint arXiv:2401.05605](#). 644 645 646

Hongyu Li, Liang Ding, Meng Fang, and Dacheng  
Tao. 2024a. Revisiting catastrophic forgetting  
in large language model tuning. [arXiv preprint](#)  
[arXiv:2406.04836](#). 647 648 649 650

Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen  
Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan  
Li, Zhengwei Tao, Xinyu Wang, and 1 others. 2025a.  
Websailor: Navigating super-human reasoning for  
web agent. [arXiv preprint arXiv:2507.02592](#). 651 652 653 654 655

Shilong Li, Xingyuan Bu, Wenjie Wang, Jiaheng Liu,  
Jun Dong, Haoyang He, Hao Lu, Haozhe Zhang,  
Chenchen Jing, Zhen Li, and 1 others. 2025b.  
Mm-browsecomp: A comprehensive benchmark  
for multimodal browsing agents. [arXiv preprint](#)  
[arXiv:2508.13186](#). 656 657 658 659 660 661

Yangning Li, Yinghui Li, Xinyu Wang, Yong Jiang,  
Zhen Zhang, Xinran Zheng, Hui Wang, Hai-Tao  
Zheng, Pengjun Xie, Philip S. Yu, Fei Huang, and  
Jingren Zhou. 2024b. Benchmarking multimodal  
retrieval augmented generation with dynamic vqa  
dataset and self-adaptive planning agent. 662 663 664 665 666 667

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,  
Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi  
Deng, Chenyu Zhang, Chong Ruan, and 1 others.  
2024. Deepseek-v3 technical report. [arXiv preprint](#)  
[arXiv:2412.19437](#). 668 669 670 671 672

673	Kartik Narayan, Yang Xu, Tian Cao, Kavya Nerella,	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	729
674	Vishal M Patel, Navid Shiee, Peter Gräsch, Chao Jia,	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	730
675	Yinfei Yang, and Zhe Gan. 2025. Deepmmsearch-	Gao, Chengen Huang, Chenxu Lv, and 1 others.	731
676	r1: Empowering multimodal llms in multimodal web	2025. Qwen3 technical report. <a href="#">arXiv preprint</a>	732
677	search. <a href="#">arXiv preprint arXiv:2510.12801</a> .	<a href="#">arXiv:2505.09388</a> .	733
678	OpenAI. 2025. Openai deep research system card.	Xiaohan Yu, Zhihan Yang, and Chong Chen. 2025.	734
679	<a href="#">OpenAI Blog</a> .	Unveiling the potential of multimodal retrieval aug-	735
680	Perplexity. 2025. Perplexity deep research system card.	mented generation with planning. <a href="#">arXiv preprint</a>	736
681	Rulin Shao, Rui Qiao, Varsha Kishore, Niklas Muen-	<a href="#">arXiv:2501.15470</a> .	737
682	nighoff, Xi Victoria Lin, Daniela Rus, Bryan	Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai,	738
683	Kian Hsiang Low, Sewon Min, Wen-tau Yih,	Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025.	739
684	Pang Wei Koh, and 1 others. 2025. Reasonir: Train-	<a href="#">Deepresearcher: Scaling deep research via reinforce-</a>	740
685	ing retrievers for reasoning tasks. <a href="#">arXiv preprint</a>	<a href="#">ment learning in real-world environments</a> . <a href="#">Preprint</a> ,	741
686	<a href="#">arXiv:2504.20595</a> .	<a href="#">arXiv:2504.03160</a> .	742
687	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	<b>A Implementation Details</b>	743
688	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	<b>A.1 RAG</b>	744
689	Zhang, YK Li, Yang Wu, and 1 others. 2024.	For the RAG baseline, we adopt the prompt de-	745
690	Deepseekmath: Pushing the limits of mathematical	sign and processing pipeline proposed in ( <a href="#">Wu et al.</a> ,	746
691	reasoning in open language models. <a href="#">arXiv preprint</a>	<a href="#">2025b</a> ). Specifically, we first perform image re-	747
692	<a href="#">arXiv:2402.03300</a> .	trieval using the Serper API, followed by query	748
693	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	refinement based on the image search context. The	749
694	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	refined query is then used to conduct a subsequent	750
695	Lin, and Chuan Wu. 2024. Hybridflow: A flexible	text search, and the final response is generated us-	751
696	and efficient rlhf framework. <a href="#">arXiv preprint arXiv:</a>	ing the combined information. It is worth noting	752
697	<a href="#">2409.19256</a> .	that our implementation differs slightly from ( <a href="#">Wu</a>	753
698	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	<a href="#">et al., 2025b</a> ) in that we do not utilize the Jina API	754
699	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	to fetch the full webpage content.	755
700	Rong Wen. 2025. R1-searcher: Incentivizing the	<b>A.2 OmniSearch</b>	756
701	search capability in llms via reinforcement learning.	For the OmniSearch baseline, we leverage the	757
702	<a href="#">arXiv preprint arXiv:2503.05592</a> .	publicly available implementation <sup>4</sup> and adapt the	758
703	Duong T Tran, Trung-Kien Tran, Manfred Hauswirth,	search tool interface to match our experimental	759
704	and Danh Le Phuoc. 2025. Reasonvqa: A multi-hop	setup. Apart from this minor adjustment, we retain	760
705	reasoning benchmark with structural knowledge for	the original workflow and prompt design to ensure	761
706	visual question answering. In <a href="#">Proceedings of the</a>	a fair comparison.	762
707	<a href="#">IEEE/CVF International Conference on Computer</a>	<b>A.3 CogPlanner</b>	763
708	<a href="#">Vision</a> , pages 18793–18803.	For CogPlanner, we develop a multi-agent plan-	764
709	Liang Wang, Nan Yang, Xiaolong Huang, Binxing	ning framework built upon the llama-index li-	765
710	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,	brary <sup>5</sup> . This implementation integrates dynamic	766
711	and Furu Wei. 2022. Text embeddings by weakly-	query reformulation and retrieval strategy selec-	767
712	supervised contrastive pre-training. <a href="#">arXiv preprint</a>	tion to facilitate efficient multimodal information	768
713	<a href="#">arXiv:2212.03533</a> .	synthesis.	769
714	Jason Wei, Zhiqing Sun, Spencer Papay, Scott McK-	<b>A.4 MMSearch-R1</b>	770
715	inney, Jeffrey Han, Isa Fulford, Hyung Won Chung,	We adopt two implementation strategies. First, we	771
716	Alex Tachard Passos, William Fedus, and Amelia	build upon the publicly available codes <sup>6</sup> and sub-	772
717	Glaese. 2025. Browsecomp: A simple yet challeng-	stitute the original image and text search compo-	773
718	ing benchmark for browsing agents. <a href="#">arXiv preprint</a>	nents with our custom-built tools. However, under	774
719	<a href="#">arXiv:2504.12516</a> .		
720	Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin,		
721	Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun		
722	Xi, Gang Fu, Yong Jiang, and 1 others. 2025a. Web-		
723	dancer: Towards autonomous information seeking		
724	agency. <a href="#">arXiv preprint arXiv:2505.22648</a> .		
725	Jinming Wu, Zihao Deng, Wei Li, Yiding Liu, Bo You,		
726	Bo Li, Zejun Ma, and Ziwei Liu. 2025b. Mmsearch-		
727	r1: Incentivizing llms to search. <a href="#">arXiv preprint</a>		
728	<a href="#">arXiv:2506.20670</a> .		

<sup>4</sup><https://github.com/Alibaba-NLP/OmniSearch>

<sup>5</sup>[https://github.com/run-llama/llama\\_index](https://github.com/run-llama/llama_index)

<sup>6</sup><https://github.com/EvolvingLMs-Lab/multimodal-search-r1>

775 this setting, the model’s tool invocation frequency  
776 rapidly diminishes to nearly zero. To ensure a fair  
777 and stable evaluation, we additionally employ the  
778 released model checkpoint<sup>7</sup> for our experiments.

## 779 B Prompts

### Multimodal RAG Rollout

You are an expert in information seeking and reasoning. You will be given a question with a image. You need to collect information for the question step by step.

Follow these instructions carefully:

1. If you need external knowledge, call search tools.
2. Enclose your entire reasoning process within <think> ... </think> tags.
3. If you find no further external knowledge needed, stop the search process and call the answer model tool.

### Answer Reward

You are an expert evaluator. You will be given:

- A question
- Several correct (golden) answer candidates
- My provided answer

Your task:

Strictly judge whether my answer is correct compared to the golden answers.

Judgement rules:

1. The meaning of my answer **\*\*must match\*\*** one of the golden answer candidates.
2. Reject or fail to answer is wrong answer.

Output format:

<reason> The reason of judgement. <Judgement> Yes or No.

Question: {question}

Golden Answer: {cand\_ans}

My Answer: {gen}

<sup>7</sup><https://huggingface.co/lmms-lab/MMSearch-R1-7B>