

Beyond Monolithic Architectures: A Multi-Agent Search and Knowledge Optimization Framework for Agentic Search

Anonymous ACL submission

Abstract

Agentic search has emerged as a promising paradigm for complex information seeking by enabling Large Language Models (LLMs) to interleave reasoning with tool use. However, prevailing systems rely on monolithic agents that suffer from structural bottlenecks, including unconstrained reasoning outputs that inflate trajectories, sparse outcome-level rewards that complicate credit assignment, and stochastic search noise that destabilizes learning. To address these challenges, we propose **M-ASK** (Multi-Agent Search and Knowledge), a framework that explicitly decouples agentic search into two complementary roles: Search Behavior Agents, which plan and execute search actions, and Knowledge Management Agents, which aggregate, filter, and maintain a compact internal context. This decomposition allows each agent to focus on a well-defined subtask and reduces interference between search and context construction. Furthermore, to enable stable coordination, M-ASK employs turn-level rewards to provide granular supervision for both search decisions and knowledge updates. Experiments on multi-hop QA benchmarks demonstrate that M-ASK outperforms strong baselines, achieving not only superior answer accuracy but also significantly more stable training dynamics.¹

1 Introduction

The rapid evolution of Large Language Models (LLMs) has fundamentally reshaped information retrieval, driving a paradigm shift from passive keyword matching to *Agentic Search* (Shi et al., 2025). Unlike traditional retrieval systems, agentic search systems function as autonomous decision-makers capable of iterative planning, external tool querying, and information synthesis to address complex, multi-hop user needs. Represented by Search-r1 (Jin et al., 2025), recent advancements (Song

¹The source code for M-ASK is available at <https://anonymous.4open.science/r/M-ASK-36F2>.

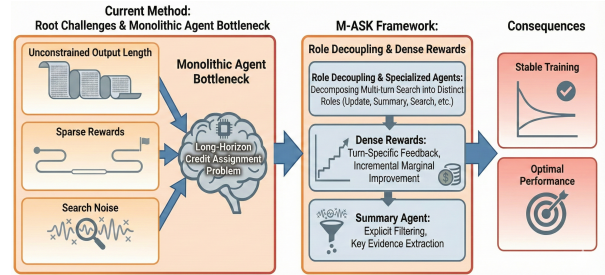


Figure 1: Challenges of current monolithic methods and the M-ASK solution. Existing agents struggle with the *long-horizon credit assignment problem* caused by unconstrained output length, sparse rewards and search noise. M-ASK addresses these bottlenecks through role decoupling and turn-level dense rewards.

et al., 2025a,b; Zheng et al., 2025) integrate reasoning directly into LLM-based multi-round search workflows. By executing multi-step search within a single response via end-to-end optimization, these approaches significantly enhance performance on complex QA tasks—capabilities that remain beyond the reach of static search paradigms (Ma et al., 2023; Ke et al., 2024).

However, training robust agents to navigate these dynamic search environments presents significant challenges. Prevailing approaches predominantly adopt a *monolithic* architecture that executes multi-round search within a single, continuous response. In this paradigm, the LLM shoulders the heavy burden of both trajectory planning and information processing at every iterative step of the generation. We argue that this monolithic design is structurally vulnerable to three intertwined obstacles: (i) **unconstrained output length**, where agents generate verbose reasoning chains that extend search horizons without necessarily increasing information density; (ii) **sparse rewards**, where feedback is typically delayed until task completion, hindering effective step-wise credit assignment; and (iii) **search noise**, where external tools such as search

066	engines introduce noise and irrelevant data into the		
067	context.		
068	As illustrated in Figure 1, these challenges are		
069	not isolated; rather, they compound to destabilize		
070	training. The interaction between extended trajec-		
071	tories and sparse, outcome-only feedback creates a		
072	severe <i>long-horizon credit assignment</i> problem: op-		
073	timization algorithms struggle to attribute the final		
074	reward to specific, distant tokens. This fragility is		
075	further exacerbated by search engine noise—when		
076	stochastic context infiltrate an already lengthy and		
077	sparsely rewarded episode, the learning signal be-		
078	comes effectively indistinguishable from variance.		
079	Consequently, monolithic agents (Jin et al., 2025)		
080	frequently suffer from suboptimal state and high		
081	training instability (Deng et al., 2025).		
082	To overcome the limitations of the monolithic		
083	method, we propose M-ASK (Multi-Agent Search		
084	and Knowledge) , a framework that fundamentally		
085	disentangles the decision-making of search from		
086	the burden of information integration. Rather than		
087	relying on a single agent, M-ASK orchestrates a		
088	collaboration between two specialized roles:		
089	1. Search Behavior Agents (including the Plan-		
090	ning, Search, and Answer Agents), which focus		
091	exclusively on trajectory planning, interacting		
092	with search engine, and generating answers;		
093	2. Knowledge Management Agents (including		
094	the Summary and Update Agents), which act as		
095	dynamic filters to prune noisy observations, up-		
096	date and maintain a concise internal knowledge		
097	state, and thus, constrain context length.		
098	Crucially, M-ASK abandons the reliance on sparse		
099	feedback in favor of turn-specific dense rewards .		
100	By jointly optimizing all agents with immediate,		
101	turn-aware supervision, we ensure that the Knowl-		
102	edge Management Agents actively stabilize the		
103	state space, empowering the Search Agents to con-		
104	duct more accurate planning. Consequently, this		
105	synergistic collaboration effectively mitigates the		
106	impact of noise and long horizons.		
107	Our contributions are summarized as follows:		
108	• We identify the compound impact of output		
109	verbosity, sparse rewards, and tool noise on		
110	agentic search, which explains why mono-		
111	lithic architectures struggle with training sta-		
112	bility (Figure 1 and Section 4.3).		
113	• We introduce M-ASK, a collaborative frame-		
114	work that decouples search behavior from		
115	knowledge management, utilizing dense, turn-		
116	specific rewards to enable joint optimization.		
		• Extensive evaluations on multi-hop QA bench-	117
		marks demonstrate that M-ASK significantly	118
		outperforms state-of-the-art baselines, deliv-	119
		ering robust gains in accuracy and markedly	120
		improving convergence stability in multi-hop	121
		search scenarios.	122
		2 Related Work	123
		2.1 From Iterative RAG to Agentic Search	124
		The paradigm of Retrieval-Augmented Genera-	125
		tion (RAG) has evolved significantly from static	126
		“retrieve-then-read” pipelines to dynamic systems.	127
		Early iterative approaches, such as IRCoT (Trivedi	128
		et al., 2023) and Self-RAG (Asai et al., 2023),	129
		introduced feedback loops where LLMs actively	130
		decide when to retrieve. Recent advancements	131
		have formalized this into Agentic Search , where	132
		models plan multi-step trajectories to solve open-	133
		ended problems using external tools. Notably,	134
		emerging <i>Deep Research</i> agents, such as DeepRe-	135
		searcher (Zheng et al., 2025), Search-o1 (Li et al.,	136
		2025), Search-r1 (Jin et al., 2025), and the R1-	137
		Searcher series (Song et al., 2025a,b), integrate	138
		retrieval directly into the reasoning chains of LLM	139
		to tackle long-horizon QA tasks.	140
		2.2 Context Management and Optimization	141
		To overcome static context limitations, re-	142
		cent works explore dynamic optimization via	143
		uncertainty-based filtering (Jimenez Gutierrez	144
		et al., 2024; Ji et al., 2025) or explicit memory	145
		agents (Yu et al., 2025; Yan et al., 2025). More	146
		aggressively, DeepNote (Wang et al., 2024) and	147
		MemSearcher (Yuan et al., 2025) treat context	148
		as an evolving state. However, DeepNote does	149
		not support joint end-to-end optimization, while	150
		MemSearcher’s monolithic design leads to coarse	151
		credit assignment due to sparse, trajectory-level	152
		rewards. M-ASK addresses these pitfalls by <i>decou-</i>	153
		<i>pling</i> search and knowledge management into two	154
		specialized agents. Through multi-agent reinforce-	155
		ment learning with turn-specific dense rewards, our	156
		approach achieves precise credit assignment.	157
		2.3 Multi-Agent Systems for Information	158
		Retrieval	159
		To overcome the limitations of monolithic agents,	160
		multi-agent systems (MAS) commonly decompose	161
		complex tasks into sub-tasks and assign them to	162
		agents with specialized roles. General frameworks,	163
		such as MetaGPT (Hong et al., 2023) and Auto-	164

Gen (Wu et al., 2024), exemplify this paradigm by enabling structured role specialization and coordination among multiple agents. Within the information retrieval domain, MindSearch (Chen et al., 2024) utilizes a graph-based planner alongside parallel web searchers for query decomposition, while MMOA-RAG (Chen et al., 2025a) allocates specialized agents for query rewriting, document selection, and answer generation. Similarly, MAO-ARAG (Chen et al., 2025b) focuses on optimizing a planner agent to balance retrieval effectiveness with efficiency. Despite their potential, most search-centric MAS rely heavily on prompt engineering or standard supervised fine-tuning (SFT), often neglecting the optimization of *interaction dynamics*. Although MMOA-RAG incorporates MAPPO (Yu et al., 2022), aims to optimize a shared global reward (e.g., the F1 score of the final answer). We argue that such sparse global feedback is suboptimal for multi-turn collaboration, as it exacerbates the *credit assignment problem*, failing to accurately evaluate intermediate steps. M-ASK addresses this limitation by introducing distinct, dense reward functions for Search and Knowledge agents, ensuring that each role is optimized for its specific contribution to the collective goal.

3 Method

We propose **M-ASK**, a multi-agent framework decoupling search planning from information integration. In this section, we first formulate the problem as a *Sequential Decentralized Partially Observable Markov Decision Process* (Bernstein et al., 2002; Oliehoek et al., 2016). Second, we detail the specifications of the specialized agents. Finally, we present the framework’s execution flow and the joint optimization process utilizing a turn-level reward mechanism.

3.1 Problem Formulation

We formulate the multi-hop QA task as a *Sequential Decentralized Partially Observable Markov Decision Process*. Unlike standard formulations where agents act simultaneously, our framework, M-ASK, operates in a turn-based manner. Specifically, at each discrete time step t , only one designated agent, denoted as π_{active} , is activated. This agent receives a partial observation of the environment and executes an action to advance the search process. To enable coordination across time steps, agents communicate indirectly by reading from and

writing to a shared, structured knowledge state.

Structured Knowledge State Agents communicate via a shared, structured *Knowledge State* \mathcal{K} . We define \mathcal{K}_t at time step t as:

$$\mathcal{K}_t = \left\{ \begin{array}{l} \text{"question"} : q, \\ \text{"thinking_trajectory"} : \mathcal{T}_t, \\ \text{"predicted_answer"} : a_t \end{array} \right\} \quad (1)$$

Let $\mathcal{T}_t = [\tau_1, \tau_2, \dots, \tau_m]$ represent the evolving reasoning chain, where each element $\tau_i = \langle q_{\text{sub}}^{(i)}, a_{\text{sub}}^{(i)} \rangle$ is a tuple consisting of a **sub-query** and its associated **sub-answer** (evidence). This trajectory \mathcal{T}_t explicitly records the step-by-step multi-hop inference process, serving as the logical derivation path required to deduce the final answer a_t from the original question q . Specifically, the **Planning Agent** initializes both the reasoning chain \mathcal{T}_0 and the answer a_0 . As the multi-round search progresses, \mathcal{T}_t is dynamically modified, and consequently, the answer a_t is iteratively updated based on the evolving trajectory. The specific implementation details are discussed in the following section. For a concrete instantiation of \mathcal{K}_t in a multi-hop scenario, please refer to Table 5 in Appendix E.

The framework shown in Figure 2 orchestrates these agents into a cohesive workflow involving initialization, iterative refinement, and optimization (see Algorithm 1 in Appendix F for the complete pseudo-code). To better understand how this workflow operates in practice, we first present the functional roles of the agents, followed by a detailed explanation of the inference and training processes.

3.2 M-ASK Architecture

M-ASK employs a team of specialized agents, categorized by their functional roles.

3.2.1 Search Behavior Agents (SBA)

Planning Agent (A_{plan}) This agent functions as the system *initializer*. Given a user query q as input, A_{plan} leverages parametric memory to generate an initial reasoning trajectory \mathcal{T}_0 and a preliminary answer a_0 , encapsulating them into the initial state $\mathcal{K}_0 \leftarrow \pi_{\text{plan}}(q)$.

Search Agent (A_{search}) Acting as the *navigator*, this agent iteratively evaluates the sufficiency of the thinking trajectory \mathcal{T}_t with respect to the query q . Its primary role is to decide between exploration and termination. Specifically, taking the question and trajectory as input, the policy outputs an action

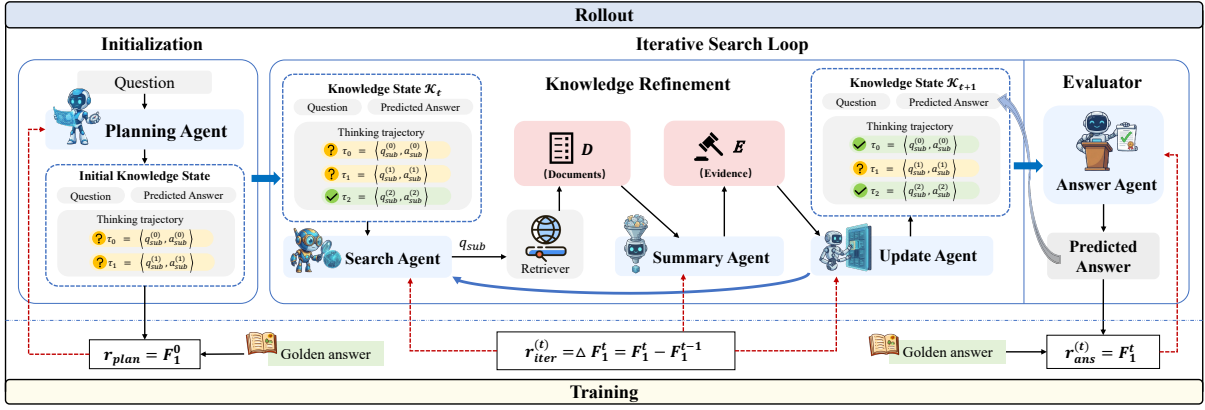


Figure 2: Overview of the M-ASK framework. (1) **Rollout**: The Planning Agent initializes the state \mathcal{K}_0 , followed by an iterative loop where Search and Knowledge Management Agents refine the trajectory. Crucially, the Answer Agent updates the prediction after each turn. (2) **Training**: A hybrid reward mechanism assigns absolute scores (F_1^0 and F_1^t) to the Planning and Answer Agents, respectively, while the collaborative agents (Search, Summary, Update) share the marginal improvement (ΔF_1^t) to incentivize step-wise refinement.

259 $Act \leftarrow \pi_{\text{search}}(q, \mathcal{T}_t)$, where $Act \in \{q'_{\text{sub}}, \langle \text{end} \rangle\}$.
 260 If expanding knowledge is necessary, it generates a
 261 specific sub-query q'_{sub} ; otherwise, it outputs $\langle \text{end} \rangle$
 262 to terminate the search loop.

263 **Answer Agent (A_{ans})** Operating as the *solver*,
 264 this agent is responsible for generating the final
 265 prediction. Conditioned on the original question
 266 q and the accumulated reasoning trajectory \mathcal{T}_t re-
 267 trieved from the knowledge state \mathcal{K} , it synthesizes
 268 a coherent final answer a , formally defined as
 269 $a \leftarrow \pi_{\text{ans}}(q, \mathcal{T}_t)$.

270 3.2.2 Knowledge Management Agents (KMA)

271 **Summary Agent (A_{sum})** This agent functions as
 272 a *filter* to distill key information. Given a sub-query
 273 q'_{sub} and a set of retrieved documents D , it extracts
 274 the pertinent evidence E while actively discarding
 275 irrelevant noise. The process is formalized as $E \leftarrow$
 276 $\pi_{\text{sum}}(q'_{\text{sub}}, D)$.

277 **Update Agent (A_{upd})** Transcending the role of
 278 a passive logger, this agent acts as a *dynamic state*
 279 *refiner*. Its primary objective is to maintain a high-
 280 density knowledge state \mathcal{K} by judiciously decid-
 281 ing between refining existing information or ap-
 282 pending new findings. Given the current state
 283 \mathcal{K}_t , a sub-query q'_{sub} , and the retrieved evidence
 284 E , the agent outputs a discrete operation op to
 285 evolve the trajectory. The action space is de-
 286 signed to balance information growth and precision:
 287 (1) $\langle \text{Update} \rangle \tau_i \langle / \text{Update} \rangle$ (**In-Place Refinement**):
 288 Targeting an existing step τ_i , this action over-
 289 writes previous hallucinations or vague informa-
 290 tion with precise evidence. (2) $\langle \text{Add} \rangle \tau_{\text{new}} \langle / \text{Add} \rangle$

(**Expansion**): This action appends a new reasoning
 291 step only when the evidence introduces a distinct,
 292 necessary logical hop. Formally, the state transition
 293 is defined as $op, \mathcal{K}_{t+1} \leftarrow \pi_{\text{upd}}(\mathcal{K}_t, q'_{\text{sub}}, E)$.
 294

295 We provide a structured summary of the func-
 296 tional roles and specifications for each agent in
 297 Table 4 (Appendix D), offering a clear view of
 298 their distinct mechanisms.

299 3.2.3 Inference Workflow

300 The inference process operates as follows:

- 301 1. **Initialization**: A_{plan} generates a "cold start"
 302 state \mathcal{K}_0 .
- 303 2. **Iteration**: The system enters a loop controlled
 304 by A_{search} .
 305 • At step t , A_{search} evaluates \mathcal{K}_t .
 306 • If A_{search} generates a query, the workflow pro-
 307 ceeds to A_{sum} (filtering) and A_{upd} (state up-
 308 date), producing \mathcal{K}_{t+1} .
- 309 3. **Termination**: The loop terminates when A_{search}
 310 outputs $\langle \text{end} \rangle$ or a maximum step limit is
 311 reached. Finally, A_{ans} is triggered to synthe-
 312 size the final answer from the latest knowledge
 313 state \mathcal{K} .

314 For a microscopic view of this workflow, includ-
 315 ing how agents dynamically correct hallucinations
 316 and refine the knowledge state, please refer to the
 317 detailed case study in Table 6 (Appendix G).

318 3.3 M-ASK Training via Turn-Level Rewards

319 We employ Independent PPO (Schulman et al.,
 320 2017) for optimization. To align the SBA and KMA
 321 groups despite their different functional roles, we

design a hybrid reward mechanism that differentiates between *state utilization* and *state refinement*.

State-Based Reward (Absolute F1) For agents responsible for generating solution outputs (A_{plan} and A_{ans}), the objective is to maximize the absolute quality of the current state.

$$r_{\text{plan}} = F_1(a_0, y), \quad r_{\text{ans}}^{(t)} = F_1(a_t, y) \quad (2)$$

Here, a_t is the answer synthesized from \mathcal{K}_t and y is the ground truth answer. Note that during training, A_{ans} acts as an evaluator at every step t .

Transition-Based Reward (Shared Incremental Gain) Crucially, the iterative phase requires tight collaboration between the Search Behavior Agent (A_{search}) and the Knowledge Management Agents ($A_{\text{sum}}, A_{\text{upd}}$). Although they belong to different functional groups, their actions are co-dependent: effective search requires precise state updates, and useful updates depend on accurate retrieval.

To enforce this **local cooperation**, we assign a **shared incremental reward** to all agents active in the loop. It is important to distinguish the execution frequency of the Answer Agent between phases. During **inference**, the Answer Agent is triggered only *once* after the search terminates to generate the final output. However, during **training**, it assumes an additional role as an *intermediate evaluator*. It is invoked at every turn t to synthesize a temporary answer a_t based on the current state \mathcal{K}_t . This mechanism allows us to utilize the answer score $r_{\text{ans}}^{(t)} = F_1(a_t, y)$ (defined in Eq. 2) to measure the immediate answer quality. The iteration reward is then defined as the marginal improvement over the previous step:

$$r_{\text{iter}}^{(t)} = r_{\text{ans}}^{(t)} - r_{\text{ans}}^{(t-1)} = F_1(a_t, y) - F_1(a_{t-1}, y) \quad (3)$$

By sharing the identical $r_{\text{iter}}^{(t)}$ signal, the Search Agent (A_{search}), Summary Agent (A_{sum}), and Update Agent (A_{upd}) are jointly incentivized to maximize the marginal information gain of each turn. This mechanism binds them into a cooperative sub-team, where the Search Agent learns to fetch necessary information and the Knowledge Agents learn to distill it efficiently. If A_{search} outputs $\langle \text{end} \rangle$, it receives a reward of 0, ensuring the team only terminates the collaboration when further search yields no positive gain.

Optimization Objective (Parameter Sharing)

To enhance sample efficiency and enable knowledge transfer across different reasoning phases, we

employ a **Parameter-Shared**² strategy. All functional agents are instantiated from a unified LLM π_θ , distinguished solely by role-specific system instructions I_{role} .

Consequently, the optimization objective aggregates the experiences from all roles. Let u_t denote the action taken given observation o_t . The shared policy π_θ is updated to maximize:

$$\mathcal{L}(\theta) = \hat{\mathbb{E}}_t \left[\min \left(\rho_t \hat{A}_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (4)$$

where $\rho_t = \frac{\pi_\theta(u_t|o_t)}{\pi_{\theta_{\text{old}}}(u_t|o_t)}$ is the probability ratio. Simultaneously, the shared Critic V_ϕ minimizes the unified value loss:

$$\mathcal{L}(\phi) = \hat{\mathbb{E}}_t \left[\|V_\phi(o_t) - R_t\|^2 \right] \quad (5)$$

Here, $R_t = \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}$ represents the cumulative discounted return, where γ is the discount factor and r denotes the role-specific rewards defined earlier.

4 Experiments

To validate the efficacy of M-ASK, we conduct comprehensive experiments to answer the following research questions: **RQ1**: Does M-ASK outperform existing monolithic and multi-agent frameworks? **RQ2**: How does M-ASK compare to the monolithic baseline (Search-r1) in terms of training stability? **RQ3**: Is it beneficial to jointly model and optimize Knowledge Management and Search Behavior? **RQ4**: Do turn-specific dense rewards provide better credit assignment than global outcome-based rewards?

4.1 Experimental Setup

Datasets We evaluate our framework on a diverse set of open-domain QA benchmarks categorized by reasoning complexity. For **Single-hop QA**, we use Natural Questions (NQ) (Kwiatkowski et al., 2019), PopQA (Mallen et al., 2022), and AmbigQA (Min et al., 2020) to test factual retrieval accuracy. For **Multi-hop QA**, to evaluate complex reasoning and trajectory planning, we employ HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020), Musique (Trivedi et al., 2022), and Bamboogle (Press et al., 2022).

²We provide a detailed justification for adopting this parameter-shared architecture in Appendix A.

Method	Single-hop QA				Multi-hop QA					Avg All
	NQ	PopQA	AmbigQA	Avg	HotpotQA	2Wiki	Musique	Bam.	Avg	
<i>Standard Baselines</i>										
LLM w/o RAG	17.53	14.93	23.53	18.66	17.76	22.58	8.58	17.14	16.52	17.44
Vanilla RAG	40.60	42.74	56.20	46.51	28.33	25.91	25.20	25.28	26.18	34.89
<i>RL-Based (Static Modular Workflow)</i>										
RRR (Ma et al., 2023)	54.60	50.46	65.41	56.82	46.21	41.52	18.27	36.59	35.65	44.72
BGM (Ke et al., 2024)	54.21	49.51	65.97	56.56	46.85	37.79	17.55	37.38	34.89	44.18
MMOA-RAG (Chen et al., 2025a)	<u>55.44</u>	<u>50.21</u>	<u>68.02</u>	<u>57.89</u>	49.21	41.66	17.26	37.20	36.33	45.57
<i>Agentic Search (Adaptive Workflow)</i>										
Adaptive RAG (Jeong et al., 2024)	36.52	35.59	45.32	39.14	42.38	39.62	<u>25.48</u>	34.85	35.58	37.11
MAO-ARAG (Chen et al., 2025b)	36.82	41.85	47.03	41.90	46.65	<u>43.96</u>	<u>22.38</u>	49.84	<u>40.85</u>	41.30
DeepNote (Wang et al., 2024)	54.03	49.80	67.57	57.13	<u>52.49</u>	36.22	22.17	<u>45.22</u>	39.03	<u>46.79</u>
Search-r1 (Jin et al., 2025)	52.57	46.98	65.25	54.93	46.87	39.03	17.97	38.69	35.64	43.91
M-ASK (Ours)	57.40	50.10	68.33	58.61	58.31	46.12	26.20	44.18	43.70	50.09
<i>Impv. vs Best</i>	+1.96	-0.36	+0.31	+0.72	+5.82	+2.16	+0.72	-5.66	+2.85	+3.30

Table 1: Main performance comparison (F1 Score) on single-hop and multi-hop QA benchmarks. The best results are **bolded** and the second best are underlined. "Impv. vs Best" denotes the performance gain (blue) or drop (gray) of M-ASK compared to the best performing baseline in each column.

Baselines We evaluate M-ASK against three categories of methods: **Standard Baselines**, **RL-Based Static Workflows**, and **Adaptive Agentic Search**. Detailed implementations and settings for these baselines are provided in Appendix B.

Implementation Details Our implementation is developed based on the official ver1 repository³. For all experiments, we use **Qwen2.5-7B-Instruct** (Team, 2024) as the backbone LLM. We utilize the English Wikipedia as our retrieval corpus, indexed via **E5** (Wang et al., 2022) for dense retrieval. We report the standard **F1 Score** as the evaluation metric.

4.2 Main Results Analysis (RQ1)

To answer **RQ1**, we evaluate M-ASK against competitive baselines across seven datasets. As shown in Table 1, M-ASK achieves the highest average F1 score (**50.09**), outperforming both monolithic and multi-agent frameworks.

Dominance on Complex QA and Generalization. M-ASK demonstrates superior capability on complex reasoning tasks, particularly on **HotpotQA**, where it surpasses the best baseline by a substantial margin of **+5.82**. This validates that our collaborative framework is far more effective at handling complex multi-hop queries than standard methods. Furthermore, M-ASK generalizes robustly to unseen out-of-domain datasets (e.g., **2Wiki**, **Musique**), consistently exceeding other Agentic baselines. This suggests M-ASK learns abstract, transferable strategies for query decompo-

sition and filtering rather than merely overfitting to training patterns.

Analysis of Bamboogle Performance. A notable exception is observed on **Bamboogle**, where M-ASK trails MAO-ARAG. We attribute this to two primary factors. First, the test set of Bamboogle is significantly smaller (only 125 queries) than other benchmarks, introducing potential statistical variance. Second, MAO-ARAG utilizes frozen pre-trained answer generators, which may answer factoid questions based on their existing parametric knowledge, whereas M-ASK is optimized to ground answers strictly in retrieved evidence. Nevertheless, M-ASK remains highly competitive, trailing the second-best method (DeepNote) by only ~ 1.0 point while outperforming all other baselines, reaffirming the robustness of our policy despite these factors.

Comparison with Monolithic and Modular Frameworks. A clear performance hierarchy is observed, further answering RQ1 regarding framework efficacy. The **monolithic** agent, Search-r1, lags significantly (43.91) due to *unconstrained context growth* and accumulated search noise. This structural limitation is addressed by **modular** frameworks like DeepNote (46.79), which achieves the second-best results by decoupling knowledge management. However, M-ASK further outperforms DeepNote (+3.30). Unlike DeepNote’s disjoint modules, M-ASK employs *end-to-end joint optimization*, ensuring that the Search and Knowledge Management agents are collaboratively updated to maximize the final reasoning reward.

³<https://github.com/volcengine/ver1>

4.3 Training Stability and Convergence (RQ2)

To investigate the training dynamics and stability of our framework compared to monolithic approaches, we visualize the training curves on the HotpotQA dataset in Figure 3. Furthermore, to rigorously quantify robustness, we conducted 10 independent training runs for both methods and recorded the rate of "model collapse"—defined as the performance score dropping to near zero and failing to recover. These statistics are summarized in Table 2.

Method	Collapse Rate (over 10 runs)		
	@ 200 Steps	@ 500 Steps	@ 1000 Steps
Search-r1	1/10 (10%)	7/10 (70%)	9/10 (90%)
M-ASK (Ours)	0/10 (0%)	0/10 (0%)	0/10 (0%)

Table 2: Training stability analysis comparing the rate of model collapse (performance degrading to ≈ 0) at different training stages across 10 independent runs.

Catastrophic Collapse in Monolithic Agents.

As shown in Table 2 and Figure 3(a), the monolithic Search-r1 agent exhibits extreme volatility. Consistent with the findings in (Deng et al., 2025), we observe that Search-r1 is prone to collapse during the mid-to-late stages of training. Specifically, while only 10% of runs failed at 200 steps, the collapse rate escalated to **90%** as training progressed to 1000 steps. We attribute this to the interaction between long horizons and sparse rewards. As the agent explores, it inevitably encounters retrieval noise. In a monolithic architecture, this noise accumulates in the context without intermediate correction. When the optimization algorithm attempts to assign credit based solely on the final outcome, the noisy gradients frequently push the policy into degenerate states from which it cannot recover.

Analysis of Response Length and Mechanism.

The key to M-ASK’s superior stability (0% collapse rate) is elucidated in Figure 3(b). While the monolithic Search-r1 suffers from "context bloating"—rapidly converging to over 1000 tokens—M-ASK maintains a remarkably low and stable average output length. This conciseness is not merely a cosmetic difference but a direct indicator of three structural advantages that stabilize training: (1) **Simplified Functionality:** Decoupling ensures atomic tasks with short generation horizons, significantly reducing optimization complexity compared to monolithic models; (2) **Precise Credit Assignment:** Short outputs allow turn-specific dense

rewards to provide immediate, high-confidence feedback, unlike global rewards that dilute over long trajectories; and (3) **Effective Noise Filtering:** Knowledge Management Agents actively prune noise to prevent error accumulation in the knowledge state.

Collectively, these factors eliminate the gradient variance caused by noisy, long contexts, and sparse reward. This structural robustness enables M-ASK to achieve the stable, monotonic convergence illustrated in Figure 3(c), even in the absence of supervised warm-starting.

4.4 Ablation Studies (RQ3 & RQ4)

To validate the contribution of individual components in M-ASK, we conduct ablation studies and observing the impact across different task complexities. The results are summarized in Table 3.

4.4.1 Impact of Collaborative Knowledge Management (RQ3)

To answer whether explicitly modeling knowledge management is beneficial, we evaluate the variant **w/o KMA**. In this setting, we remove the *Summary* and *Update* agents; the *Search Behavior Agent* interacts directly with the raw search engine output and appends full documents to the context. This setup simulates the unconstrained information flow of a standard monolithic interaction. However, unlike the fully monolithic *Search-r1*, this variant retains the underlying multi-agent paradigm and independent optimization methods, thereby strictly isolating the impact of the missing knowledge management module.

Analysis. The results show a consistent performance degradation across both single-hop (avg. $\Delta - 2.58\%$) and multi-hop (avg. $\Delta - 2.86\%$) benchmarks. This confirms that without the active filtering provided by the KMA group, the search agent is overwhelmed by noise in the retrieval results, leading to hallucinations or distracted reasoning. We note a slight anomaly on **2WikiMultiHopQA** (+0.46%), where the unfiltered model performs marginally better. However, this is outweighed by the severe drops in more complex datasets like **Musique** (−6.81%) and **HotpotQA** (−3.25%). Overall, explicit knowledge management is crucial for stabilizing the reasoning trajectory against retrieval noise.

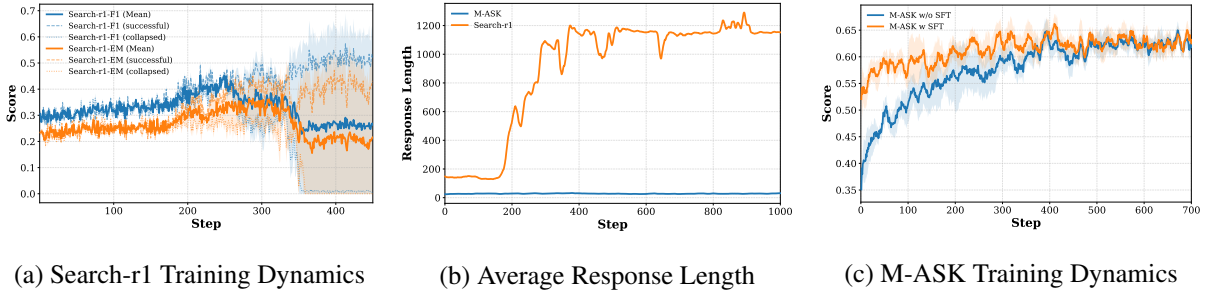


Figure 3: Training curves on HotpotQA. In (a) and (c), solid lines represent the mean across multiple runs, while shaded regions indicate the variance. (a) Search-r1 exhibits high instability and frequent mode collapse. (b) Evolution of average response length; M-ASK remains concise while Search-r1 suffers from context bloating. (c) M-ASK demonstrates stable convergence.

Method	Single-hop Datasets								Multi-hop Datasets									
	NQ		PopQA		AmbigQA		Avg.		HotpotQA		2Wiki		Musique		Bam.		Avg.	
	F1	Δ	F1	Δ	F1	Δ	F1	Δ	F1	Δ	F1	Δ	F1	Δ	F1	Δ	F1	Δ
M-ASK (Ours)	57.40	-	50.10	-	68.33	-	58.61	-	58.31	-	46.12	-	26.20	-	44.18	-	43.70	-
<i>Ablation Studies</i>																		
w/o KMA	54.93	-2.47	48.36	-1.74	64.79	-3.54	56.03	-2.58	55.06	-3.25	46.58	+0.46	19.39	-6.81	42.33	-1.85	40.84	-2.86
w/o T-L Reward	54.95	-2.45	47.12	-2.98	62.04	-6.29	54.70	-3.91	42.57	-15.74	33.27	-12.85	11.95	-14.25	26.10	-18.08	28.47	-15.23

Table 3: Ablation study. **Avg.** columns show the aggregated performance on Single-hop (3 datasets) and Multi-hop (4 datasets) benchmarks. Δ represents the performance drop (gray) or gain (blue) compared to the full model.

4.4.2 Efficacy of Turn-Specific Dense Rewards (RQ4)

To assess the necessity of our reward shaping mechanism, we evaluate **w/o T-L Reward**. In this variant, we disable the turn-level incremental rewards (Eq.2 and Eq.3) and instead optimize all agents using only the final outcome reward (Global F1) at the end of the episode, similar to the strategy used in MMOA-RAG (Chen et al., 2025a).

Analysis. This ablation reveals a critical insight into the optimization dynamics of multi-turn search. **Single-hop Resilience:** On single-hop datasets, the performance drop is moderate (avg. $\Delta - 3.91\%$) because the short trajectories typical of these tasks make the final reward sufficient for credit assignment. **Multi-hop Collapse:** In stark contrast, performance collapses on multi-hop datasets (avg. $\Delta - 15.23\%$), with the score on **HotpotQA** plummeting by over 15 points.

This disparity highlights the **turn-level credit assignment problem**. In frameworks like MMOA-RAG, all agents across different turns share an identical global outcome reward. This ambiguity obscures the specific contribution of individual steps, making it difficult to discern which action actually drove the success or failure. Consequently, agents

struggle to optimize intermediate sub-goals in long-horizon tasks. In contrast, M-ASK’s turn-specific Δ F1 reward provides immediate, dense feedback, enabling agents to accurately lock in optimal policies for each reasoning hop.

5 Conclusion

In this paper, we addressed the structural instability of monolithic agentic search caused by unconstrained contexts, sparse rewards, and search noise. We proposed **M-ASK**, a multi-agent framework that decouples search planning from knowledge management, enabling synergistic collaboration via turn-specific dense rewards.

Empirical evaluations across seven benchmarks demonstrate that M-ASK consistently outperforms state-of-the-art baselines, particularly in complex multi-hop scenarios. Crucially, our analysis reveals that decoupling these roles significantly enhances training stability compared to end-to-end RL approaches. These findings suggest that explicit role specialization and intermediate supervision are critical for scaling agentic search to more open-ended and noisy real-world environments. Future work will explore extending M-ASK to heterogeneous model architectures and broader tool-use scenarios beyond information retrieval.

615 Limitations

616 Despite the effectiveness of M-ASK, several lim-
617 itations remain. First, the collaborative multi-
618 agent workflow inherently increases the frequency
619 of LLM invocations compared to single-pass ap-
620 proaches, as each reasoning step requires discrete
621 inference calls for search, summarization, and up-
622 date modules. This sequential interaction pattern
623 inevitably raises computational costs and may pose
624 challenges for latency-sensitive applications. Sec-
625 ond, our evaluation is currently confined to textual
626 QA tasks; the framework’s generalizability to other
627 complex domains, such as code generation or mul-
628 timodal reasoning, remains to be explored. Finally,
629 the efficacy of our parameter-sharing strategy relies
630 on the inherent capacity of the backbone LLM to
631 handle diverse role instructions, and performance
632 on smaller-scale architectures warrants further in-
633 vestigation.

634 Ethics Statement

635 This work builds upon large language models to
636 facilitate multi-agent collaboration for textual ques-
637 tion answering. All datasets used in our experi-
638 ments are publicly available and do not contain
639 any personally identifiable information. Neverthe-
640 less, outputs generated by LLMs may reflect bi-
641 ases present in their pretraining data or produce
642 incorrect information. Users and practitioners
643 should exercise caution when applying the pro-
644 posed framework to real-world scenarios, espe-
645 cially in domains where misinformation or biased
646 content could cause harm. We do not foresee ad-
647 ditional ethical risks beyond those commonly associ-
648 ated with contemporary language model research.

649 References

650 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
651 Hannaneh Hajishirzi. 2023. Self-rag: Learning to
652 retrieve, generate, and critique through self-reflection.
653 *arXiv preprint arXiv:2310.11511*.

654 Daniel S Bernstein, Robert Givan, Neil Immerman, and
655 Shlomo Zilberstein. 2002. The complexity of decen-
656 tralized control of markov decision processes. *Math-*
657 *ematics of operations research*, 27(4):819–840.

658 Yiqun Chen, Hangyu Mao, Jiaxin Mao, Shiguang Wu,
659 Tianle Zhang, Bin Zhang, Wei Yang, and Hongxing
660 Chang. 2022. Ptd: Personalized training with dis-
661 tilled execution for multi-agent reinforcement learn-
662 ing. *arXiv preprint arXiv:2210.08872*.

Yiqun Chen, Lingyong Yan, Weiwei Sun, Xinyu Ma,
Yi Zhang, Shuaiqiang Wang, Dawei Yin, Yiming
Yang, and Jiaxin Mao. 2025a. Improving retrieval-
augmented generation through multi-agent reinforce-
ment learning. *arXiv preprint arXiv:2501.15228*.

Yiqun Chen, Erhan Zhang, Lingyong Yan, Shuaiqiang
Wang, Jizhou Huang, Dawei Yin, and Jiaxin Mao.
2025b. Mao-arag: Multi-agent orchestration for
adaptive retrieval-augmented generation. *arXiv*
preprint arXiv:2508.01005.

Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning
Liu, Wenwei Zhang, Kai Chen, and Feng Zhao. 2024.
Mindsearch: Mimicking human minds elicits deep ai
searcher. *arXiv preprint arXiv:2407.20183*.

Wenlong Deng, Yushu Li, Boying Gong, Yi Ren, Chris-
tos Thrampoulidis, and Xiaoxiao Li. 2025. On
grpo collapse in search-r1: The lazy likelihood-
displacement death spiral. *arXiv preprint*
arXiv:2512.04220.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,
and Akiko Aizawa. 2020. Constructing a multi-hop
qa dataset for comprehensive evaluation of reasoning
steps. *arXiv preprint arXiv:2011.01060*.

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu
Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang,
Zili Wang, Steven Ka Shing Yau, Zijuan Lin, and
1 others. 2023. Metagpt: Meta programming for a
multi-agent collaborative framework. In *The Twelfth*
International Conference on Learning Representa-
tions.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam
Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,
Akila Welihinda, Alan Hayes, Alec Radford, and 1
others. 2024. Gpt-4o system card. *arXiv preprint*
arXiv:2410.21276.

Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju
Hwang, and Jong C Park. 2024. Adaptive-rag: Learn-
ing to adapt retrieval-augmented large language mod-
els through question complexity. *arXiv preprint*
arXiv:2403.14403.

Yuelyu Ji, Rui Meng, Zhuochun Li, and Daqing He.
2025. Memory-aware and uncertainty-guided re-
trieval for multi-hop question answering. *arXiv*
preprint arXiv:2503.23095.

Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michi-
hiro Yasunaga, and Yu Su. 2024. Hipporag: Neu-
robiologically inspired long-term memory for large
language models. *Advances in Neural Information*
Processing Systems, 37:59532–59569.

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.

717	Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang,	John Schulman, Filip Wolski, Prafulla Dhariwal,	770
718	Qiaozhu Mei, and Michael Bendersky. 2024. Bridg-	Alec Radford, and Oleg Klimov. 2017. Proxi-	771
719	ing the preference gap between retrievers and llms.	mal policy optimization algorithms. <i>arXiv preprint</i>	772
720	<i>arXiv preprint arXiv:2401.06954</i> .	<i>arXiv:1707.06347</i> .	773
721	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Zhengliang Shi, Yiqun Chen, Haitao Li, Weiwei Sun,	774
722	field, Michael Collins, Ankur Parikh, Chris Alberti,	Shiyu Ni, Yougang Lyu, Run-Ze Fan, Bowen Jin,	775
723	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kent-	Yixuan Weng, Minjun Zhu, Qiuqie Xie, Xinyu Guo,	776
724	on Lee, and 1 others. 2019. Natural questions: a	Qu Yang, Jiayi Wu, Jujia Zhao, Xiaqiang Tang, Xin-	777
725	benchmark for question answering research. <i>Trans-</i>	bei Ma, Cunxiang Wang, Jiaxin Mao, and 7 others.	778
726	<i>actions of the Association for Computational Linguis-</i>	2025. <i>Deep research: A systematic survey</i> . <i>Preprint</i> ,	779
727	<i>tics</i> , 7:453–466.	<i>arXiv:2512.02038</i> .	780
728	Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang,	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	781
729	Yujia Zhou, Yutao Zhu, Peitian Zhang, and	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	782
730	Zhicheng Dou. 2025. Search-o1: Agentic search-	Rong Wen. 2025a. R1-searcher: Incentivizing the	783
731	enhanced large reasoning models. <i>arXiv preprint</i>	search capability in llms via reinforcement learning.	784
732	<i>arXiv:2501.05366</i> .	<i>arXiv preprint arXiv:2503.05592</i> .	785
733	Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI	Huatong Song, Jinhao Jiang, Wenqing Tian, Zhipeng	786
734	Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent	Chen, Yuhuan Wu, Jiahao Zhao, Yingqian Min,	787
735	actor-critic for mixed cooperative-competitive envi-	Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025b.	788
736	ronments. <i>Advances in neural information process-</i>	R1-searcher++: Incentivizing the dynamic knowl-	789
737	<i>ing systems</i> , 30.	edge acquisition of llms via reinforcement learning.	790
738	Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao,	<i>arXiv preprint arXiv:2505.17005</i> .	791
739	and Nan Duan. 2023. Query rewriting for retrieval-	Qwen Team. 2024. Qwen2 technical report. <i>arXiv</i>	792
740	augmented large language models. <i>arXiv preprint</i>	<i>preprint arXiv:2412.15115</i> .	793
741	<i>arXiv:2305.14283</i> .	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	794
742	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi	and Ashish Sabharwal. 2022. Musique: Multi-	795
743	Das, Hannaneh Hajishirzi, and Daniel Khashabi.	hop questions via single-hop question composition.	796
744	2022. When not to trust language models: Invest-	<i>Transactions of the Association for Computational</i>	797
745	igating effectiveness and limitations of paramet-	<i>Linguistics</i> , 10:539–554.	798
746	ric and non-parametric memories. <i>arXiv preprint</i>	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	799
747	<i>arXiv:2212.10511</i> , 7.	and Ashish Sabharwal. 2023. Interleaving retrieval	800
748	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and	with chain-of-thought reasoning for knowledge-	801
749	Luke Zettlemoyer. 2020. Ambigqa: Answering	intensive multi-step questions. In <i>Proceedings of</i>	802
750	ambiguous open-domain questions. <i>arXiv preprint</i>	<i>the 61st annual meeting of the association for com-</i>	803
751	<i>arXiv:2004.10645</i> .	<i>putational linguistics (volume 1: long papers)</i> , pages	804
752	Frans A Oliehoek, Christopher Amato, and 1 oth-	10014–10037.	805
753	ers. 2016. <i>A concise introduction to decentralized</i>	Liang Wang, Nan Yang, Xiaolong Huang, Binxing	806
754	<i>POMDPs</i> , volume 1. Springer.	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,	807
755	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt,	and Furu Wei. 2022. Text embeddings by weakly-	808
756	Noah A Smith, and Mike Lewis. 2022. Measuring	supervised contrastive pre-training. <i>arXiv preprint</i>	809
757	and narrowing the compositionality gap in language	<i>arXiv:2212.03533</i> .	810
758	models. <i>arXiv preprint arXiv:2210.03350</i> .	Ruobing Wang, Qingfei Zhao, Yukun Yan, Daren Zha,	811
759	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	Yuxuan Chen, Shi Yu, Zhenghao Liu, Yixuan Wang,	812
760	pher D Manning, Stefano Ermon, and Chelsea Finn.	Shuo Wang, Xu Han, and 1 others. 2024. Deep-	813
761	2023. Direct preference optimization: Your language	note: Note-centric deep retrieval-augmented genera-	814
762	model is secretly a reward model. <i>Advances in neural</i>	tion. <i>arXiv preprint arXiv:2410.08821</i> .	815
763	<i>information processing systems</i> , 36:53728–53741.	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	816
764	Tabish Rashid, Mikayel Samvelyan, Chris-	Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,	817
765	tian Schroeder De Witt, Gregory Farquhar,	Shaokun Zhang, Jiale Liu, and 1 others. 2024. Au-	818
766	Jakob Foerster, and Shimon Whiteson. 2020.	togen: Enabling next-gen llm applications via multi-	819
767	Monotonic value function factorisation for deep	agent conversations. In <i>First Conference on Lan-</i>	820
768	multi-agent reinforcement learning. <i>Journal of</i>	<i>guage Modeling</i> .	821
769	<i>Machine Learning Research</i> , 21(178):1–51.	Sikuan Yan, Xiufeng Yang, Zuchao Huang, Ercong Nie,	822
		Zifeng Ding, Zonggen Li, Xiaowen Ma, Kristian Ker-	823
		sting, Jeff Z Pan, Hinrich Schütze, and 1 others. 2025.	824
		Memory-r1: Enhancing large language model agents	825

826 to manage and utilize memories via reinforcement
827 learning. *arXiv preprint arXiv:2508.19828*.

828 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-
829 gio, William W Cohen, Ruslan Salakhutdinov, and
830 Christopher D Manning. 2018. Hotpotqa: A dataset
831 for diverse, explainable multi-hop question answer-
832 ing. *arXiv preprint arXiv:1809.09600*.

833 Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao,
834 Yu Wang, Alexandre Bayen, and Yi Wu. 2022. The
835 surprising effectiveness of ppo in cooperative multi-
836 agent games. *Advances in Neural Information Pro-
837 cessing Systems*, 35:24611–24624.

838 Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie
839 Chen, Weinan Dai, Qiyong Yu, Ya-Qin Zhang, Wei-
840 Ying Ma, Jingjing Liu, Mingxuan Wang, and 1 others.
841 2025. Memagent: Reshaping long-context llm with
842 multi-conv rl-based memory agent. *arXiv preprint
843 arXiv:2507.02259*.

844 Qianhao Yuan, Jie Lou, Zichao Li, Jiawei Chen, Yaojie
845 Lu, Hongyu Lin, Le Sun, Debing Zhang, and Xianpei
846 Han. 2025. Memsearcher: Training llms to reason,
847 search and manage memory via end-to-end reinforce-
848 ment learning. *arXiv preprint arXiv:2511.02805*.

849 Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai,
850 Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025.
851 Deepresearcher: Scaling deep research via reinforce-
852 ment learning in real-world environments. *arXiv
853 preprint arXiv:2504.03160*.

854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903

A Rationale for Parameter Sharing Strategy

In the M-ASK framework, we employ a parameter-sharing strategy where a single Large Language Model (LLM) backbone π_θ serves as the underlying policy for all functional agents (i.e., A_{plan} , A_{search} , A_{sum} , A_{upd} , A_{ans}), distinguished solely by role-specific system instructions. We adopt this design based on three critical considerations:

1. Theoretical Foundation in MARL Parameter sharing is a well-established and effective paradigm in Multi-Agent Reinforcement Learning (MARL) (Rashid et al., 2020; Lowe et al., 2017; Yu et al., 2022; Chen et al., 2022). In cooperative settings, sharing parameters allows agents to learn a unified representation of the state space, which often leads to faster convergence and improved training stability compared to maintaining independent policies. By optimizing a single set of parameters on the aggregated experiences of all roles, the model can efficiently generalize across the diverse phases of the search and reasoning process.

2. Computational and Storage Efficiency Unlike traditional RL agents based on small Multi-Layer Perceptrons (MLPs), agents in our framework are initialized with LLMs containing billions of parameters. Maintaining independent policy networks for five distinct agents would result in a linear increase in memory consumption ($\mathcal{O}(N)$), rendering the training process computationally prohibitive and difficult to deploy. Parameter sharing reduces the storage requirement to $\mathcal{O}(1)$, significantly lowering the barrier for training and inference. This efficiency is paramount for complex RAG systems like M-ASK, allowing us to allocate resources toward longer context windows rather than redundant model weights.

3. Inherent Multi-Task Capability of LLMs LLMs inherently possess strong multi-task capabilities, enabling them to perform distinct tasks based on contextual instructions (prompts) without modifying their internal weights. In M-ASK, different agents (e.g., the *Search Agent* deciding on queries vs. the *Summary Agent* extracting evidence) share fundamental reasoning competencies, such as reading comprehension and logical deduction. Parameter sharing leverages this synergy: skills learned in the *summarization* task can implicitly enhance

the *answer generation* capability via shared representations. By conditioning the shared π_θ on role-specific prompts I_{role} , we effectively project the model’s general capabilities into specific functional subspaces, achieving role specialization without architectural redundancy.

B Implementation Details of Baselines

We compare M-ASK against three categories of methods:

- **Standard Baselines:** *LLM w/o RAG* (closed-book), *Vanilla RAG* (standard retrieve-then-generate).
- **RL-Based (Static Modular Workflow):** *RRR* (Ma et al., 2023) (RL-based query reformulation), and *BGM* (Ke et al., 2024) (RL-based document selection), and *MMOA-RAG* (Chen et al., 2025a) (multi-agent RL).
- **Agentic Search (Adaptive Workflow):** *Adaptive RAG* (Jeong et al., 2024) (adaptive workflow), *MAO-ARAG* (Chen et al., 2025b) (planner-executors optimization), *DeepNote* (Wang et al., 2024) (knowledge-management) and *Search-r1* (Jin et al., 2025) (monolithic agent optimized via RL reasoning training).

To ensure a fair comparison, all baselines and our M-ASK method are unified under the same experimental setting. We utilize Qwen2.5-7B-Instruct as the backbone model. Specifically, for components within the baselines that do not require training, we employ the pre-trained version of Qwen2.5-7B-Instruct. For components requiring training, we fine-tune them based on the Qwen2.5-7B-Instruct initialization. Detailed implementation notes for each baseline are provided below:

Standard Baselines

- **LLM w/o RAG:** This represents the closed-book setting where the LLM directly generates answers based on its internal parametric knowledge without accessing external corpora.
- **Vanilla RAG:** A standard retrieve-then-generate pipeline. It retrieves the top-5 relevant documents based on the query and concatenates them with the prompt to generate the answer using the pre-trained LLM.

904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950

Static Modular Workflow (RL-Based)

- **RRR** (Ma et al., 2023): This method employs the PPO algorithm to end-to-end train a query rewriter module. In our reproduction, to prevent the performance bottleneck potentially caused by a frozen generator, we also fine-tuned the answer generation module, ensuring it is well-trained alongside the rewriter.
- **BGM** (Ke et al., 2024): BGM uses PPO to train a document selection module. Although the original implementation may utilize a frozen generator, we fine-tuned the answer generation module in our reproduction to align with the robust setting of RRR and maximize overall performance.
- **MMOA-RAG** (Chen et al., 2025a): This framework utilizes Multi-Agent PPO (MAPPO) (Yu et al., 2022) to train three distinct agents—a query rewriter, a document selector, and an answer generator—guided by a shared final reward mechanism. We adopted Qwen2.5-7B-Instruct as the backbone model for this multi-agent framework, strictly following the original logic while maintaining consistency with our unified setting.

Agentic Search (Adaptive Workflow)

- **Adaptive RAG** (Jeong et al., 2024): This method trains a classifier to dynamically determine which of three distinct retrieval workflows (Directly Answer, Simple RAG, or Iterative RAG) should handle a given query. In our implementation, we fine-tuned Qwen2.5-7B-Instruct to serve as this classifier, while utilizing the pre-trained version of Qwen2.5-7B-Instruct for executing the subsequent static workflows.
- **MAO-ARAG** (Chen et al., 2025b): MAO-ARAG features a hierarchical multi-agent architecture composed of a planner and multiple executors, sharing a modular structure similar to our M-ASK. The original work utilizes PPO with a shared final reward to train specifically the planner agent. Consistent with our unified setting, the planner module is a fine-tuned Qwen2.5-7B-Instruct, while the executor modules employ the pre-trained version of the same model.

- **DeepNote** (Wang et al., 2024): DeepNote focuses on optimizing knowledge management within QA tasks through Direct Preference Optimization (DPO) (Rafailov et al., 2023). In our reproduction, we upgraded the data generation model from gpt-4o-mini (Hurst et al., 2024) (used in the original paper) to gpt-4o (Hurst et al., 2024) to construct higher-quality DPO training data.
- **Search-r1** (Jin et al., 2025): Unlike modular designs, Search-r1 adopts a monolithic architecture where multi-turn search and reasoning processes occur within a single response generation. It is optimized via RL using outcome-based rewards. For this baseline, we obtained the experimental results using the code provided in the author’s official open-source repository⁴.

C Prompt for Different Agents

C.1 Prompt for Planning Agent

The following is the full system prompt used for decomposing queries. The [EXAMPLE_PROMPT] placeholder represents the few-shot examples injected at runtime.

System Prompt

```
You are a knowledgeable reasoning assistant. Your task is to decompose a given question into multiple sub-questions based purely on your own knowledge (do not use external tools), and answer each sub-question based on your own knowledge. If you do not know the answer, write 'unkown' exactly.
```

```
Output all results strictly in the following format for easy parsing:  
<q1>sub-question text</q1>  
<a1>answer text</a1>  
<q2>sub-question text</q2>  
<a2>answer text</a2>
```

```
Continue numbering sequentially until done. At the end, output your final conclusion answer to the original question in the tag <predicted_answer>...</predicted_answer>, and inside this tag follow the style of: [EXAMPLE_PROMPT]
```

User Input

```
Query: {query}
```

⁴<https://github.com/PeterGriffinJin/Search-R1>

C.2 Prompt for Search Agent

System Prompt

You are a search query generation assistant.
Your input contains:

- The original question
- A `thinking_trajectory` (list of sub-questions and their sub-answers)
- A `search_history` (list of past search queries and their summaries)

Note:

- The `search_history` represents queries YOU have already searched in previous turns.
- You **MUST** avoid generating queries that are identical to or semantically similar to any query in `search_history`.
- Only generate a new query if it provides new, distinct information that has not yet been searched.

Your task:

1. If any sub-question in `thinking_trajectory` has an answer 'unkown' or unclear, generate a specific, searchable query to find the missing information.
2. If you believe additional sub-questions are needed to answer the original question, generate a specific searchable query for that.
3. If no further search is needed, output `<end>`.

Output only one of the following two formats:

```
<search>Your query here</search>
<end>
```

Do not include any extra text, explanations, or other tags.

User Input

Original question: {question}

Thinking trajectory:
{trajectory_text}

Search history:
{history_text}

Output only one of the following two formats:

```
<search>Your query here</search>
<end>
```

Do not include any extra text, explanations, or other tags.

C.3 Prompt for Summary Agent

System Prompt

You are an evidence extraction assistant.
You will be given:

- A search query
- Several candidate documents related to that query

Your task:

1. From the provided documents, extract the key evidence that directly answers the query.
2. If no key evidence exists in the documents, output 'No useful information'.

Output exactly in one of the two formats:
<evidence>your key evidence here</evidence>
<evidence>No useful information</evidence>

Do not include explanations, reasoning steps, or any text outside the <evidence> tags.

User Input

Query: {query}
Candidate Documents: {docs_text}

C.4 Prompt for Update Agent

System Prompt

You are an update detection assistant.
You are given:

- An original question and its `thinking_trajectory` (sub-questions and their answers)
- The `predicted_answer` derived from `thinking_trajectory`
- A new question and its evidence

Your task:

1. If the new question and evidence can replace one existing sub-question-answer pair in `thinking_trajectory`, output `<Update>ti</Update>` where `ti` is the key.
2. If the new question and evidence contains information not present in `thinking_trajectory`, output `<Add>t(n+1)</Add>` where `n` is current number of `t`'s.
3. Output strictly one of these formats and no other text.

User Input

Original question: {orig_question}
Thinking trajectory: {trajectory_text}
New query: {new_query}
Evidence: {new_evidence}

C.5 Prompt for Answer Agent

System Prompt

You are a final answer generation assistant.
You are given:
- An original question
- A `thinking_trajectory` (list of sub-questions and their answers)

Your task:
1. Use the provided `thinking_trajectory`'s answers to determine the final answer to the original question.
2. You must only output in the following strict format as shown in [EXAMPLE_PROMPT]

IMPORTANT: Output must be:
<predicted_answer>your answer
here</predicted_answer>
No extra text, no reasoning.

User Input

Question: {question}
Thinking trajectory: {trajectory_text}

D Detailed Agent Specifications

In this section, we provide comprehensive specifications for the multi-agent architecture employed in M-ASK. Due to space constraints in the main text, we present the granular details of the Search Behavior Agents (SBA) and Knowledge Management Agents (KMA) here. Table 4 formally defines the specific roles, input-output interfaces, and action spaces for each agent, elucidating their operational logic within the iterative retrieval process.

E Example of Structured Knowledge State

To facilitate a better understanding of the data structure defined in Eq. 1, we provide a concrete example of the Structured Knowledge State \mathcal{K}_t .

Consider a multi-hop question q : "Who is the current CEO of the company that developed ChatGPT?". Table 5 illustrates the state \mathcal{K}_t , where the agent has successfully decomposed the query and retrieved the necessary evidence chain.

F Joint Training Algorithm Details

In this section, we provide the comprehensive pseudo-code for the M-ASK joint training process, complementing the methodology described in Section 3.3. Algorithm 1 details the execution flow

of the **Parameter-Shared** strategy, where a unified policy π_θ is employed across all functional roles (Planning, Search, Summary, Update, and Answer).

The procedure consists of two distinct phases:

- **Data Collection Phase:** Agents interact sequentially to solve the task. Notably, the Answer Agent serves as an intermediate evaluator at each step t , calculating the marginal performance gain $\Delta F1$ to provide dense, turn-specific supervision.
- **Unified Optimization Phase:** Experiences from all roles are stored in a mixed replay buffer \mathcal{B} . The shared model parameters θ and critic ϕ are then jointly updated via PPO, maximizing the objective function derived from these heterogeneous interaction trajectories.

G Detailed Case Study: Structured Knowledge State Evolution

The trajectory presented in Table 6 offers a microscopic view of how M-ASK addresses the structural limitations of monolithic agents. The execution process exhibits high rationality and ingenuity in three critical aspects:

1. Correction of Parametric Hallucinations via "Targeting". A pervasive challenge in agentic search is the "cold start" problem, where models must plan without initial observations. In the **Init** phase, the Planning Agent relies on its parametric memory and incorrectly predicts "Guadalajara." Crucially, M-ASK does not treat this initial plan as ground truth but as a *search target*. In **Turn 1**, the Update Agent identifies a conflict between the parametric "Guadalajara" and the retrieved non-parametric evidence ("Puebla"). Instead of attempting to reconcile the two or hedging, the agent prioritizes external evidence, executing an <Update> action to overwrite the hallucination. This demonstrates a robust **error-correction mechanism** that prevents initial errors from propagating through the reasoning chain.

2. Deep Verification of Temporal Constraints. The transition to **Turn 2** highlights the Search Agent's ability to handle complex constraints. The query contains a specific temporal condition ("1943"). While a naive retriever might stop after finding the entity "Puebla," the Search Agent

Module	Agent	Role	Description & Action Space	Input / Output Formulation
SBA	Planning	Initializer	Leverages parametric memory to generate an initial reasoning trajectory \mathcal{T}_0 and a preliminary answer a_0 , encapsulating them into the initial state.	In: Query q Out: $\mathcal{K}_0 \leftarrow \pi_{\text{plan}}(q)$
	Search	Navigator	Evaluates information sufficiency to decide between exploration and termination. 1) Generate sub-query q'_{sub} to expand knowledge. 2) Output <end> to terminate the loop.	In: Query q , Trajectory \mathcal{T}_t Out: $Act \leftarrow \pi_{\text{search}}(\mathcal{K}_t)$, where $Act \in \{q'_{\text{sub}}, \text{<end>}\}$
	Answer	Solver	Synthesizes the final answer prediction conditioned on the original question and the accumulated reasoning trajectory.	In: Query q , Trajectory \mathcal{T}_t Out: $a \leftarrow \pi_{\text{ans}}(q, \mathcal{T}_t)$
KMA	Summary	Filter	Distills pertinent evidence from retrieved documents while actively discarding irrelevant noise to ensure information density.	In: Sub-query q'_{sub} , Documents D Out: $E \leftarrow \pi_{\text{sum}}(q'_{\text{sub}}, D)$
	Update	Dynamic Refiner	Judiciously decides how to integrate evidence to evolve the trajectory: • <Update> τ_i </Update> (In-Place Refinement): Overwrites hallucinations/vague steps. • <Add> τ_{new} </Add> (Expansion): Appends necessary logical hops.	In: State \mathcal{K}_t , sub-query q'_{sub} , Evidence E Out: $op, \mathcal{K}_{t+1} \leftarrow \pi_{\text{upd}}(\mathcal{K}_t, q'_{\text{sub}}, E)$

Table 4: Comprehensive specifications of the agents within the M-ASK framework, detailing their roles, functional mechanisms, and input/output formalizations.

Key	Value (Content)
"question" (q)	Who is the current CEO of the company that developed ChatGPT?
"thinking_trajectory" (\mathcal{T}_t)	(τ_1) : $\langle q_{\text{sub}}^{(1)} \rangle$: Which company developed ChatGPT? $\langle a_{\text{sub}}^{(1)} \rangle$: OpenAI developed ChatGPT, an AI chatbot launched in November 2022. (τ_2) : $\langle q_{\text{sub}}^{(2)} \rangle$: Who is the current CEO of OpenAI? $\langle a_{\text{sub}}^{(2)} \rangle$: Sam Altman is the CEO of OpenAI.
"predicted_answer" (a_t)	Sam Altman

Table 5: A example of the Knowledge State \mathcal{K}_t . The trajectory \mathcal{T}_t contains a sequence of sub-query and sub-answer pairs that logically support the final answer.

actively verifies the historical context ("significant event... in 1943"). The subsequent <Add> t_2 </Add> action is highly rational: the agent recognizes that the "professionalization of the league" is distinct background information that supports the answer's validity. By appending this as a new step (t_2), the system constructs a logical evidence chain: *Entity Existence* (t_1) + *Historical Context* (t_2).

3. Active Prevention of Context Bloating. The most significant display of ingenuity occurs in **Turn 3**. The agent retrieves specific confirmation linking "Puebla" to the "Athletic Club." In standard monolithic frameworks (e.g., Search-r1), this new retrieval would typically be appended to the context, causing the input length to grow linearly with search steps. In contrast, M-ASK employs **In-**

Place Refinement. The Update Agent recognizes that this information is a more precise version of t_1 . It triggers <Update> t_1 to refine the query and answer (adding "specifically Puebla") rather than creating a redundant t_3 . This mechanism maintains a concise state space (\mathcal{K}_3), effectively solving the "unconstrained output length" challenge and ensuring the final reasoning is performed on a high-density, noise-free context.

Algorithm 1: M-ASK Joint Training (Parameter Shared)

Input: Data \mathcal{D} , Unified Model π_θ , Critic V_ϕ , Role Instructions $I = \{I_{\text{plan}}, I_{\text{search}}, \dots\}$, Batch B

Output: Optimized Parameters θ^*, ϕ^*

```
1 Initialize  $\theta, \phi$ ; Initialize unified replay buffer  $\mathcal{B}$ ;  
2 for iteration  $i = 1, \dots, M$  do  
  // 1. Data Collection Phase  
3  while  $|\mathcal{B}| < B$  do  
4    Sample  $(q, y) \sim \mathcal{D}$ ;  
    // Phase I: Initialization  
5     $\mathcal{K}_0 \leftarrow \pi_\theta(q; I_{\text{plan}})$ ;  $a_0 \leftarrow \mathcal{K}_0.\text{predicted\_answer}$ ;  
6     $S_{\text{prev}} \leftarrow \text{F1}(a_0, y)$ ;  
7    Add transition  $(I_{\text{plan}}, q, \mathcal{K}_0, r = S_{\text{prev}})$  to  $\mathcal{B}$ ;  
    // Phase II: Iterative Collaboration  
8     $t \leftarrow 0$ ;  
9    while  $t < T_{\text{max}}$  do  
10     // Step 1: SBA Decision  
     Action  $\leftarrow \pi_\theta(\mathcal{K}_t; I_{\text{search}})$ ; // Prompt as Searcher  
11     if Action == <end> then  
12       Add transition  $(I_{\text{search}}, \mathcal{K}_t, \text{<end>, } r = 0)$  to  $\mathcal{B}$ ;  
13       break;  
14     end  
15      $q'_{\text{sub}} \leftarrow \text{Action}$ ;  
     // Step 2: KMA Execution  
16      $D \leftarrow \text{SearchEngine}(q'_{\text{sub}})$ ;  
17      $E \leftarrow \pi_\theta(q'_{\text{sub}}, D; I_{\text{sum}})$ ; // as Summarizer  
18      $op, \mathcal{K}_{\text{next}} \leftarrow \pi_\theta(\mathcal{K}_t, q'_{\text{sub}}, E; I_{\text{upd}})$ ; // as Updater  
     // Step 3: Evaluation  
19      $a_{\text{next}} \leftarrow \pi_\theta(\mathcal{K}_{\text{next}}; I_{\text{ans}})$ ; // as Answerer  
20      $S_{\text{curr}} \leftarrow \text{F1}(a_{\text{next}}, y)$ ;  
21      $\Delta\text{F1} \leftarrow S_{\text{curr}} - S_{\text{prev}}$ ;  
     // Step 4: Store Mixed Experiences  
22     Add  $(I_{\text{search}}, \mathcal{K}_t, q'_{\text{sub}}, r = \Delta\text{F1})$  to  $\mathcal{B}$ ;  
23     Add  $(I_{\text{sum}}, q'_{\text{sub}}, D, E, r = \Delta\text{F1})$  to  $\mathcal{B}$ ;  
24     Add  $(I_{\text{upd}}, \mathcal{K}_t, E, op, r = \Delta\text{F1})$  to  $\mathcal{B}$ ;  
25     Add  $(I_{\text{ans}}, \mathcal{K}_{\text{next}}, a_{\text{next}}, r = S_{\text{curr}})$  to  $\mathcal{B}$ ;  
26      $\mathcal{K}_{t+1} \leftarrow \mathcal{K}_{\text{next}}$ ;  $S_{\text{prev}} \leftarrow S_{\text{curr}}$ ;  $t \leftarrow t + 1$ ;  
27   end  
28 end  
  // 2. Unified Optimization Phase  
29  Compute GAE on unified buffer  $\mathcal{B}$  using shared Critic  $V_\phi$ ;  
30  for epoch  $k = 1, \dots, K$  do  
31    Sample mixed mini-batches from  $\mathcal{B}$ ;  
32    Update  $\theta, \phi$  via PPO maximizing joint objective  $\mathcal{L}(\theta)$ ;  
33  end  
34  Clear buffer  $\mathcal{B}$ ;  
35 end
```

Turn	Role	Full Action Output & Structured Knowledge State (\mathcal{K}_t)
Target Question: What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? Golden Answer: Puebla		
Init	Planning Agent	$\langle q \rangle$ What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? $\langle /q \rangle$ $\langle a \rangle$ Guadalajara $\langle /a \rangle$ $\langle \text{predicted_answer} \rangle$ Guadalajara $\langle / \text{predicted_answer} \rangle$
		State \mathcal{K}_0 (Initial Structure): "question": What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? "thinking_trajectory": [$[\tau_1] \langle q_{sub}^{(1)} \rangle$: What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? $\langle a_{sub}^{(1)} \rangle$: Guadalajara (Parametric Hallucination)] "predicted_answer": Guadalajara
Turn 1	Search Agent	$\langle \text{search} \rangle$ What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? $\langle / \text{search} \rangle$ <i>Search Engine</i> [External Tool] Retrieved 5 document fragments (~100 words each) containing keyword matches for "Mexican city", "Athletic Club", "Primera Fuerza", "1943".
	Summary Agent	$\langle \text{evidence} \rangle$ Puebla, which has a city of Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza. $\langle / \text{evidence} \rangle$
	Update Agent	$\langle \text{Update} \rangle \tau_1 \langle / \text{Update} \rangle$ (Corrects Hallucination)
	Answer Agent	$\langle \text{predicted_answer} \rangle$ Puebla $\langle / \text{predicted_answer} \rangle$
		State \mathcal{K}_1 (Corrected Answer): "question": What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? "thinking_trajectory": [$[\tau_1] \langle q_{sub}^{(1)} \rangle$: What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? $\langle a_{sub}^{(1)} \rangle$: Puebla, which has a city of Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza.] "predicted_answer": Puebla (Updated via retrieval)
Turn 2	Search Agent	$\langle \text{search} \rangle$ What is a significant event or fact about Puebla A.C.'s football team's participation in the Primera Fuerza in 1943? $\langle / \text{search} \rangle$
	<i>Search Engine</i>	[External Tool] Retrieved 5 document fragments (~100 words each) related to "Puebla A.C.", "Primera Fuerza participation", "1943 events".
	Summary Agent	$\langle \text{evidence} \rangle$ Puebla A.C.'s football team played in the Primera Fuerza prior to the professionalization and development of the Primera División de México in 1943. $\langle / \text{evidence} \rangle$
	Update Agent	$\langle \text{Add} \rangle \tau_2 \langle / \text{Add} \rangle$ (Appends Context)
	Answer Agent	$\langle \text{predicted_answer} \rangle$ Puebla $\langle / \text{predicted_answer} \rangle$
		State \mathcal{K}_2 (Expanded Context): "question": What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? "thinking_trajectory": [$[\tau_1] \langle q_{sub}^{(1)} \rangle$: What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? $\langle a_{sub}^{(1)} \rangle$: Puebla, which has a city of Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza. , $[\tau_2] \langle q_{sub}^{(2)} \rangle$: What is a significant event or fact about Puebla A.C.'s football team's participation in the Primera Fuerza in 1943? $\langle a_{sub}^{(2)} \rangle$: Puebla A.C.'s football team played in the Primera Fuerza prior to the professionalization and development of the Primera División de México in 1943.] "predicted_answer": Puebla
Turn 3	Search Agent	$\langle \text{search} \rangle$ What Mexican city, specifically Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza in 1943? $\langle / \text{search} \rangle$
	<i>Search Engine</i>	[External Tool] Retrieved 5 document fragments verifying "Puebla", "Puebla A.C.", "Primera Fuerza 1943".
	Summary Agent	$\langle \text{evidence} \rangle$ Puebla, specifically Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza in 1943. $\langle / \text{evidence} \rangle$
	Update Agent	$\langle \text{Update} \rangle \tau_1 \langle / \text{Update} \rangle$ (Refines Existing Step)
	Answer Agent	$\langle \text{predicted_answer} \rangle$ Puebla $\langle / \text{predicted_answer} \rangle$
		State \mathcal{K}_3 (Refined Trajectory): "question": What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943? "thinking_trajectory": [$[\tau_1] \langle q_{sub}^{(1)} \rangle$: What Mexican city, specifically Puebla , had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza in 1943? $\langle a_{sub}^{(1)} \rangle$: Puebla, specifically Puebla, had an Athletic Club, Puebla A.C., whose football team played in the Primera Fuerza in 1943. , $[\tau_2] \langle q_{sub}^{(2)} \rangle$: What is a significant event or fact about Puebla A.C.'s football team's participation in the Primera Fuerza in 1943? $\langle a_{sub}^{(2)} \rangle$: Puebla A.C.'s football team played in the Primera Fuerza prior to the professionalization and development of the Primera División de México in 1943.] "predicted_answer": Puebla
Turn 4	Search Agent	$\langle \text{end} \rangle$ (Termination Triggered)
Final Prediction: Puebla (✓ Matches Golden Answer)		

Table 6: Full execution log for the query: "What Mexican city had an Athletic Club whose football team played in the Primera Fuerza in 1943?". The **State** rows explicitly show the three components of \mathcal{K}_t . The trajectory steps are indexed as $[\tau_i]$, allowing the Update Agent to target specific steps (e.g., $\langle \text{Update} \rangle \tau_1 \langle / \text{Update} \rangle$).