

UR²: Unify RAG and Reasoning through Reinforcement Learning

Weitao Li^{1,2,*}, Boran Xiang³, Xiaolong Wang^{1,2}, Jingyi Ren^{1,2}, Ante Wang²,
Zhinan Gou³, Weizhi Ma^{2,†}, Yang Liu^{1,2,†}

¹ Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua University, Beijing, China

² Institute for AI Industry Research (AIR), Tsinghua University, Beijing, China

³ School of Management Science & Information Engineering, Hebei University of Economics and Business, Hebei, China

Abstract

Large Language Models (LLMs) have shown strong capabilities through two complementary paradigms: Retrieval-Augmented Generation (RAG) for knowledge grounding and Reinforcement Learning from Verifiable Rewards (RLVR) for complex reasoning. However, existing attempts to unify these paradigms remain narrow in scope, typically limited to open-domain QA with fixed retrieval settings, which constrains generalization to broader domains. To address this limitation, we propose UR² (Unified RAG and Reasoning), a general reinforcement learning framework that dynamically coordinates retrieval and reasoning. UR² introduces two key designs: a difficulty-aware curriculum that selectively invokes retrieval only for challenging instances, and a hybrid knowledge access strategy that combines domain-specific offline corpora with on-the-fly LLM-generated summaries. Together, these components mitigate the imbalance between retrieval and reasoning and improve robustness to noisy information. Experiments on open-domain QA, MMLU-Pro, medical, and mathematical reasoning tasks show that UR², built on Qwen-2.5-3/7B and LLaMA-3.1-8B, consistently outperforms existing RAG and RL baselines, and achieves performance comparable to GPT-4o-mini and GPT-4.1-mini on several benchmarks. We will release all code, models, and data.

1 Introduction

Large Language Models (LLMs) have achieved remarkable performance across diverse tasks by incorporating external knowledge (Retrieval-Augmented Generation, RAG) (Lewis et al., 2020; Borgeaud et al., 2022; Izacard et al., 2022) and optimizing reasoning through reinforcement learning with verifiable rewards (RLVR) (Guo et al., 2025a). RAG methods enable LLMs to access external knowledge, while RLVR shows strong gains

on mathematical and logical reasoning (Zeng et al., 2025; Chen et al., 2025). Motivated by these successes, recent work has begun to integrate retrieval and reasoning: for example, Search-o1 (Li et al., 2025b) embeds an agentic RAG workflow into the LLM’s chain-of-thought, and RAG-Gym (Xiong et al., 2025) proposes a unified RL-based training framework for RAG agents. Similarly, RAG-RL methods which learn to invoke retrieval through RL—such as R1-Searcher (Song et al., 2025a) and Search-R1 (Jin et al., 2025) use RLVR to train models on *when* and *what* to retrieve during reasoning, improving performance in open-domain QA.

Despite recent progress, RAG-RL frameworks remain limited in scope. Most methods focus narrowly on open-domain QA, with retrieval tied to fixed reasoning steps or static knowledge sources like Wikipedia. However, paradigms that work well on open-domain QA often fail to transfer to broader domains. **Two key limitations persist:** (1) **over-reliance on retrieval**, where models issue ill-posed or trivial queries to offload reasoning and weaken internal inference; (2) **collapse to pure reasoning**, where retrieval is avoided due to noisy documents and ineffective integration, causing the model to degenerate into standalone Chain-of-Thought reasoning. For instance, R1-Searcher and Search-R1 assume access to Wikipedia, ill-suited for tasks requiring specialized information. While methods like DeepResearcher attempt training in real web environments, they face inefficiencies due to the noisy and unstructured nature of online data (Zheng et al., 2025). Other methods like ZeroSearch (Sun et al., 2025a), use LLM-generated corpora to simulate retrieval, avoiding API costs but risking hallucination and loss of real-world complexity.

To address the limitations of existing RAG-RL approaches such as static retrieval, limited domain generalization, and poor robustness in noisy environments—we propose a general and adaptive framework, UR² (Unified RAG and Reasoning),

*Email: liwt23@mails.tsinghua.edu.cn

which uses RL to dynamically coordinate retrieval and reasoning. Unlike prior methods that rely solely on static corpora (e.g., Wikipedia) or simulate retrieval with synthetic content, UR² combines both: it leverages task-specific offline corpora for accurate grounding, augmented with on-the-fly LLM-generated summaries for efficiency and generalization. To address the **imbalance** between retrieval and reasoning in prior methods, we design a difficulty-aware curriculum that adaptively controls when to trigger retrieval during training. Specifically, retrieval is used only for hard instances, encouraging the model to rely on internal reasoning when possible and to learn retrieval strategies only when necessary. This reduces retrieval overhead, improves query quality on challenging questions, and preserves reasoning capability across tasks.

We train UR² on Qwen-2.5-3B/7B-Instruct (Yang et al., 2024) and LLaMA-3.1-8B-Instruct (Dubey et al., 2024) across MMLU-Pro, Medicine, Math, and open-domain QA. During training, these models spontaneously develop key cognitive behaviors: self-verification through retrieval, intermediate reasoning validation, and hypothesis revision based on external evidence. UR² outperforms previous state-of-the-art (SOTA) methods by **5.8%** (7B) and **19.0%** (3B) on average, with peak gains of **9.5%** and **29.6%**. Notably, our 7B model matches GPT-4o-mini and GPT-4.1-mini¹, and generalizes well across domains and model architectures.

Our main contributions are summarized as follows:

- We propose the first unified retrieval-reasoning RL framework that adapts to diverse tasks beyond open-domain QA, representing an important milestone for AI systems combining parametric and external knowledge.
- We design a unified data representation and training scheme bridging retrieval and reasoning, with difficulty-aware curricula and LLM-based summarization of retrieved evidence for accurate grounding and broad generalization, providing a solution to the imbalance between retrieval and reasoning.
- Comprehensive experiments demonstrate that UR² surpasses advanced RAG and RL meth-

ods without expert demonstrations and generalizes robustly across domains.

2 Related Work

2.1 Retrieval-Augmented Generation

RAG enhances LLMs by incorporating external information to reduce hallucinations (Gao et al., 2023). Early RAG methods concatenate retrieved documents with input prompts (Lewis et al., 2020; Izacard et al., 2022; Borgeaud et al., 2022). Subsequent approaches have evolved in multiple directions: advanced RAG methods incorporate sophisticated retrieval and re-ranking strategies (Gao et al., 2023; Peng et al., 2024); post-hoc verification methods address hallucinations by retrieving documents based on generated responses (Li et al., 2024; Sun et al., 2024); and Graph-based RAG methods integrate knowledge graphs for multihop reasoning (Edge et al., 2024; Hu et al., 2025b; Peng et al., 2024). Recent RAG-RL frameworks have explored retrieval integration during training via real-time or synthetic retrieval (Zheng et al., 2025; Sun et al., 2025a). However, these approaches remain constrained by static retrieval strategies, limited domain generalization, and **inability to dynamically coordinate retrieval with reasoning** across diverse task types.

2.2 Reinforcement Learning for Retrieval-Enhanced Reasoning

RL has emerged as a key technique for significantly improving LLM capabilities, evolving from early policy gradient methods such as REINFORCE (Williams, 1992) to more advanced algorithms like PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024). Recent methods, including ARENA (Ren et al., 2025), Search-R1 (Jin et al., 2025), and R1-Searcher (Song et al., 2025a), demonstrate that RL enables LLMs to effectively learn multi-step reasoning and retrieval strategies without requiring human feedback. These works collectively highlight a clear shift from fixed retrieval heuristics to learned, RL-driven retrieval policies, which form the foundation for our unified framework, with retrieval becoming **increasingly parameterized** rather than merely prompt-guided.

3 Method

We propose UR², a unified framework that integrates retrieval-based grounding with explicit step-by-step reasoning through reinforcement learning.

¹<https://chat.openai.com/>

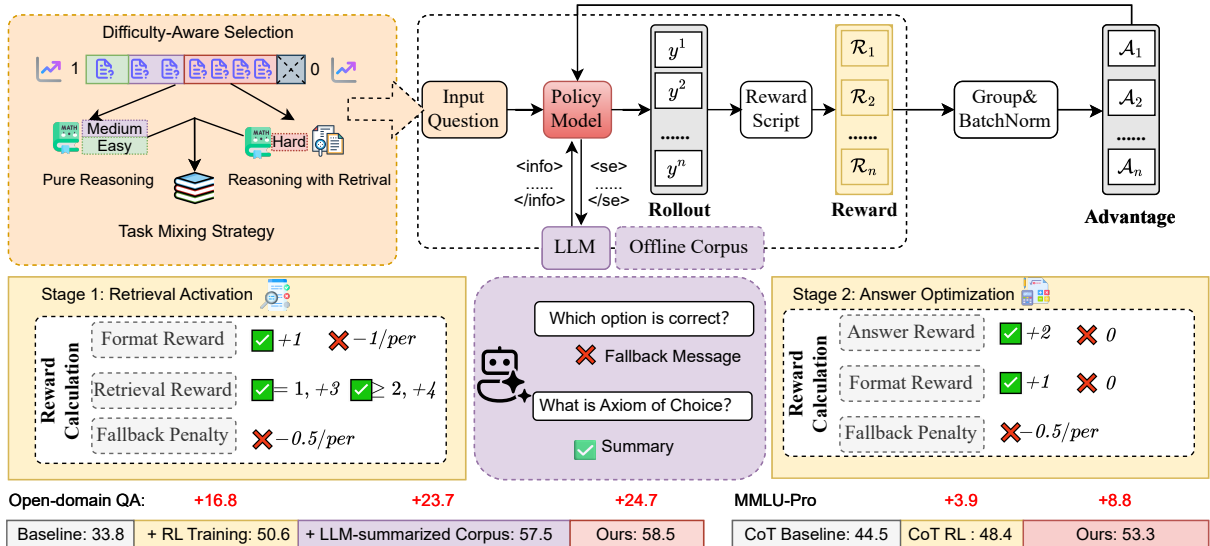


Figure 1: Overview of the UR² training pipeline. The top illustrates LLM-summarized retrieval corpus, difficulty-aware curriculum design and a two-stage reward design for retrieval activation and answer optimization. The bottom horizontal bars indicate: (1) Open-domain QA: ablation results under F1 score; and (2) MMLU-Pro: comparison with baselines using EM score.

When extending retrieval-augmented generation to reasoning-intensive domains (e.g., mathematical problem solving and domain-specific QA), we observe a fundamental challenge: *an imbalance between retrieval and reasoning*. Specifically, models may either over-rely on retrieval, offloading reasoning to external evidence, or collapse to pure chain-of-thought reasoning while failing to utilize external knowledge when necessary.

UR² is designed to address this imbalance in a principled and systematic manner. Rather than treating retrieval as an always-on module or a purely auxiliary component, our framework explicitly coordinates *when* to retrieve, *how* retrieved information is presented, and *which training signals encourage balanced behavior*. To this end, UR² combines (i) on-the-fly LLM-based summarization that transforms retrieved evidence into compact, reasoning-compatible representations, and (ii) a difficulty-aware curriculum that selectively activates retrieval based on task hardness and knowledge demands. Together, these components enable UR² to generalize across reasoning-intensive tasks while avoiding both retrieval overuse and reasoning degradation.

3.1 Imbalance Between Retrieval and Reasoning

When extending retrieval-augmented generation to reasoning-intensive tasks in our experiments, we observe a systematic imbalance between retrieval and reasoning behaviors: 1) **Over-reliance**

on retrieval. In a controlled setting, we remove constraints on retrieval usage and allow the model to freely issue retrieval queries during RL training. We observe that the model frequently generates ill-posed or trivial queries (e.g., “which option is correct”), effectively offloading reasoning to retrieval and weakening its internal reasoning process. This behavior leads to unstable training and degraded downstream performance, indicating that unconstrained retrieval is harmful in reasoning-centric tasks. 2) **Collapse to pure reasoning**. Conversely, under a Search-R1–style training setup, where retrieval summaries are removed and the model is trained with raw or noisy retrieved documents, training becomes unstable and retrieval is gradually abandoned. In this case, the model degenerates to pure chain-of-thought reasoning, failing to leverage external knowledge even for knowledge-intensive questions. This phenomenon highlights that retrieval noise directly interferes with reasoning and discourages effective retrieval usage.

These observations indicate that effective RAG for reasoning requires explicit mechanisms to regulate retrieval activation and its interaction with the reasoning process.

3.2 Preventing Over-Reliance on Reasoning via LLM-Summarized Retrieval

When retrieval is exposed to raw or noisy documents, models trained with RL tend to abandon retrieval altogether and collapse to pure chain-of-thought reasoning. To prevent this failure mode,

You are solving a multiple-choice question. Analyze each option step by step and select the best choice. If you're uncertain about any fact, you may issue a **search query** like this: `<se>` a concise query (under 20 words) `</se>`

- You may issue **multiple queries** during your reasoning.
- Each query should focus on **only one specific fact or concept** and **Avoid combining multiple facts in a single query.**

[Examples omitted here]

- You may use **up to four queries** in total — use them wisely.

When documents are returned in the format: `<info>` ... (search results here) `</info>`, integrate the retrieved information into your reasoning to refine your analysis and reach a well-supported conclusion. Finally, give your answer in this format: the correct answer is: A, B, C, D, etc.

Figure 2: Instruction prompt used to guide retrieval-augmented reasoning in UR². See Appendix D.2 for details.

UR² introduces an LLM-summarized retrieval mechanism that converts retrieved content into compact and reasoning-compatible evidence.

UR² accesses domain-relevant knowledge sources, including domain-specific resources (e.g., curated medical references or encyclopedic content). Instead of directly feeding retrieved documents to the model, retrieved content is transformed into concise, reasoning-compatible representations using an LLM, avoiding the introduction of misleading or low-quality evidence.

Beyond noise reduction, the summarization mechanism provides a substantial computational benefit. In online retrieval settings, retrieved documents average 11,346 tokens, which the summarizer compresses to 761 tokens—a **14.91×** compression ratio (see Appendix B.6). Given the quadratic complexity of the attention mechanism, this compression trades a small sequential summarization cost for a much larger reduction in reasoning-time computation, often making the overall pipeline faster despite the additional summarization step. This design ensures that retrieval remains a usable and beneficial capability during optimization, preventing degeneration into pure reasoning even in knowledge-intensive or noisy environments.

3.3 Preventing Over-Reliance on Retrieval via Difficulty-Aware Curriculum

Conversely, unconstrained retrieval usage during training can cause models to offload reasoning to external search, leading to shallow reasoning behaviors and ill-posed retrieval queries. To avoid this failure mode, UR² employs a difficulty-aware curriculum that explicitly regulates when retrieval is used during training.

Training Data Selection We categorize training instances by difficulty to control retrieval exposure. For each question, we perform 20 roll-outs using Qwen-2.5-7B-Instruct and compute the average performance score (s), following (Song et al., 2025a; Huang et al., 2025; Sun et al., 2025b). Based on s , questions are grouped into Easy ($0.8 \leq s \leq 1.0$), Medium ($0.5 \leq s < 0.8$), and Hard ($0.2 \leq s < 0.5$) levels. Following prior work (Yu et al., 2025; Guo et al., 2025b), extremely difficult samples ($s < 0.2$) are filtered, as they hinder stable learning. We adopt a sampling ratio of 7:2:1 for hard, medium, and easy questions, prioritizing challenging instances while retaining sufficient direct reasoning cases.

Task Mixing Strategy Based on difficulty, we selectively assign retrieval-augmented or pure reasoning training formats. For mathematical reasoning tasks, only hard problems invoke retrieval-augmented prompting (Figure 2), while easy and medium problems rely on direct step-by-step reasoning. In contrast, open-domain QA consistently uses retrieval due to its inherent knowledge dependency.

By associating retrieval usage with task difficulty, this curriculum prevents retrieval from becoming a default behavior. The model learns to rely on external knowledge only when necessary, preserving intrinsic reasoning ability while still benefiting from retrieval on knowledge-intensive problems. More detailed experimental configurations are provided in the Section 4 and Appendix C.2.

3.4 Two-Stage Optimization for UR²

Given the limited tool invocation capabilities of base models, especially in reasoning-integrated scenarios, we design a two-stage optimization framework to systematically develop retrieval skills and

reasoning proficiency. We train UR² using REINFORCE++ (Guo et al., 2025a), a streamlined variant of PPO. To prevent overfitting to retrieved content, we adopt retrieval masking (Song et al., 2025a; Jin et al., 2025). Our implementation is based on the REINFORCE++-baseline provided by OpenRLHF (Hu et al., 2024). See Appendix C.2 and Section 4.4 for detailed implementation.

RAG-based Rollout UR² enables the model to issue retrieval queries during reasoning rather than pre-retrieving all information upfront. As illustrated in Figure 2, the prompting mechanism enforces key principles: queries target single facts grounded in external knowledge, retrieval occurs when needed during the reasoning process, and strict format constraints using special tokens demarcate retrieval actions.

This design allows the model to strategically leverage external knowledge by learning *when* to retrieve and *what* to query for purposeful and grounded reasoning.

Stage 1: Retrieval Capability Activation We use UR² with Qwen-2.5-7B-Instruct on mathematical and open-domain QA tasks as an example. In Stage 1, the model trains on mathematical problems requiring retrieval calls in the specified format (Figure 2). The objective is not answer accuracy, but to enforce correct usage of the retrieval tool and promote retrieval-invoking behavior. This specialized training runs for only 10 steps. Further details on task setup and extensions to other models are provided in Appendix C.7.

The total reward is:

$$R_{i,\text{stage1}} = R_{i,\text{format}} + R_{i,\text{retrieval}} - P_{i,\text{fallback}} \quad (1)$$

where (1) **Format Reward**: +1 for fully compliant output; -1 per violation (e.g., malformed tags, overlength queries, missing retrieval, or illegal tokens); (2) **Retrieval Reward**: +3 for one valid query, +4 for two or more; (3) **Fallback Penalty**: -0.5 per fallback fault.

This stage equips the model with retrieval capabilities and promotes effective integration of retrieved information during generation.

Stage 2: Answer Quality Optimization Building on Stage 1, we incorporate correctness feedback to refine generation quality while preserving retrieval behaviors. The updated reward function is:

$$R_{i,\text{stage2}} = R_{i,\text{answer}} + R_{i,\text{format}} - P_{i,\text{fallback}} \quad (2)$$

where (1) **Answer Reward**: +2 for correct answers, 0 for incorrect; (2) **Format Reward**: +1 for fully valid format; 0 otherwise; (3) **Fallback Penalty**: -0.5 per fallback fault.

By decoupling retrieval skill acquisition (Stage 1) from generation optimization (Stage 2), we ensure stable convergence and interpretable credit assignment across complex reasoning trajectories.

4 Experimental Settings

4.1 Training Datasets

We build a unified training set covering math (SimpleZoo-RL (Zeng et al., 2025)), open-domain QA (R1-Searcher (Song et al., 2025a)), and multi-choice medical QA (MedQA (Jin et al., 2021)). For MMLU-Pro (Wang et al., 2024a) domains (philosophy, history, economics), we generate synthetic questions via Qwen-3-32B². After deduplication and data selection using Qwen-2.5-7B-Instruct on 20 rollouts per question, we obtain 3K samples each for math, open-domain QA, and MedQA, and 2K samples for each MMLU-Pro domain. Details are in Appendix C.2.

4.2 Evaluation Benchmarks

We evaluate generalization across four task families: (1) **Math Reasoning**: MATH500 (in-domain) (Hendrycks et al., 2021), Minerva (OOD) (Lewkowycz et al., 2022); metric: LLM-as-a-judge. (2) **Medical QA**: MedQA (5-choice, in) (Jin et al., 2021), MMLU-Pro Medical (M-Med, OOD); metric: EM. (3) **MMLU-Pro**: Philosophy, History, Economics (in), Law (OOD); metric: EM. (4) **Open-Domain QA**: HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020) (in); Bamboogle (Press et al., 2023), MusiQue (Trivedi et al., 2022) (OOD); metrics: F1 and LLM-as-a-judge.

4.3 Baselines

We compare UR² to: (1) **Vanilla Methods**: Chain-of-Thought (Kojima et al., 2022), Standard RAG (Borgeaud et al., 2022; Izacard et al., 2022) (top- $k=10$). (2) **Advanced RAG Methods**: Search-o1 (Li et al., 2025b), Self-Ask (Press et al., 2023), and RAT (Wang et al., 2024b), which combine reasoning with retrieval using prompt. (3) **CoT-RL Methods**: R1-like methods including Open-Reasoner-Zero (Hu et al., 2025a), SimpleRL-Zoo (Zeng et al., 2025), and

²<https://huggingface.co/Qwen/Qwen3-32B>

Method	MMLU-Pro				Medicine			Math			
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	Avg	MedQA [†]	M-Med [‡]	Avg	Math500 [†]	Minerva [‡]	Avg
GPT-4o-mini											
CoT	56.7	<u>53.1</u>	<u>70.4</u>	38.2	<u>54.5</u>	71.4	67.0	69.2	<u>78.0</u>	65.6	71.8
Standard RAG	<u>57.0</u>	52.3	68.6	<u>36.2</u>	53.5	70.6	64.2	67.4	77.1	68.4	<u>72.8</u>
<i>Advanced RAG Methods</i>											
Self-Ask	56.3	48.5	67.8	31.2	51.0	72.4	<u>68.0</u>	70.2	62.9	45.2	54.1
RAT	57.5	55.3	73.0	34.2	55.0	<u>74.4</u>	70.6	72.5	77.5	64.2	70.9
Search-o1	53.5	55.3	69.8	35.4	53.5	75.2	66.6	<u>70.9</u>	78.6	<u>68.3</u>	73.5
Qwen-2.5-7B											
CoT	42.3	45.7	63.4	26.6	44.5	57.2	52.0	54.6	76.6	59.4	68.0
Standard RAG	44.6	41.1	57.8	26.0	42.4	54.2	53.2	53.7	73.8	54.6	64.2
<i>CoT-RL Methods</i>											
General Reasoner	47.9	44.2	65.9	30.4	47.1	58.4	54.4	56.4	76.6	62.1	68.8
Open-Reasoner-Zero	50.0	46.6	67.5	<u>34.2</u>	49.6	61.6	58.8	60.2	<u>80.7</u>	58.8	<u>69.8</u>
SimpleRL-Zoo	35.7	36.9	55.2	25.4	38.3	57.2	51.0	54.1	77.1	50.7	63.9
<i>Our Implementations</i>											
Vanilla RL	<u>52.2</u>	43.5	64.0	33.8	48.4	64.2	57.4	60.8	78.2	59.4	68.8
UR ² (qw Summary)	51.2	<u>49.7</u>	<u>67.8</u>	32.6	<u>50.3</u>	<u>67.2</u>	<u>60.8</u>	<u>64.0</u>	80.2	59.4	<u>69.8</u>
UR ² (GPT Summary)	53.2	53.0	72.2	35.0	53.3	69.6	62.8	66.2	80.9	<u>61.0</u>	71.0
Qwen-2.5-3B											
CoT	33.6	32.3	48.8	20.6	33.8	39.4	36.8	38.1	63.6	39.9	51.8
Standard RAG	37.8	36.5	51.4	23.2	37.1	45.6	40.0	42.8	65.3	40.8	53.1
<i>Our Implementations</i>											
Vanilla RL	40.7	34.7	55.0	24.6	38.7	51.8	47.6	49.7	<u>68.0</u>	43.9	56.0
UR ² (qw summary)	<u>43.0</u>	<u>42.6</u>	<u>60.6</u>	<u>26.6</u>	<u>43.2</u>	<u>55.2</u>	<u>54.0</u>	<u>54.6</u>	68.4	47.4	57.9
UR ² (GPT summary)	47.8	49.3	63.9	30.0	47.8	59.8	56.8	58.3	67.8	<u>45.0</u>	<u>57.2</u>
LLaMA-3.1-8B											
CoT	37.8	<u>40.9</u>	53.4	29.0	40.3	59.6	52.6	56.1	48.4	34.4	41.4
Standard RAG	43.6	33.9	51.0	26.6	38.8	56.4	53.2	54.8	45.0	31.4	38.2
<i>Our Implementations</i>											
Vanilla RL	<u>44.6</u>	36.9	53.0	26.4	40.2	<u>66.8</u>	<u>57.4</u>	<u>62.1</u>	45.5	43.4	<u>44.4</u>
UR ² (qw summary)	43.6	41.2	<u>56.8</u>	28.0	<u>42.4</u>	65.8	56.2	61.0	48.6	38.4	43.5
UR ² (GPT summary)	48.3	38.6	58.0	<u>28.8</u>	43.4	68.6	58.4	63.5	54.5	<u>39.0</u>	46.8

Table 1: Performance on reasoning and math tasks. We report EM scores (in %) on MMLU-Pro and MedQA, and LLM-as-a-judge scores (in %) on math benchmarks. † = in-domain, ‡ = out-of-domain. Best results are **bold**; second-best are underlined.

General-Reasoner (Ma et al., 2025). (4) **RAG-RL Methods**: R1-Searcher (Song et al., 2025a), R1-Searcher++ (Song et al., 2025b), Search-R1 (Jin et al., 2025), and ZeroSearch (Sun et al., 2025a). (5) **Vanilla RL**: Baseline implementation following the same training setup and datasets as UR², with RAG-RL applied to open-domain QA and CoT-RL to mathematical and multiple-choice tasks. Details are in Appendix C.3.

We use Qwen-2.5-3B-Instruct, Qwen-2.5-7B-Instruct, LLaMA-3.1-8B-Instruct, GPT-4o-mini, and GPT-4.1-mini as backbones (see Appendix C.4 for configs).

4.4 Implementation Details

Retrieval uses BGE-large-en-v1.5³ and the KILT (Petroni et al., 2021) Wikipedia corpus (100-word segments, 29M documents) following (Song et al., 2025a). Open-domain QA uses Wikipedia

abstract corpus⁴. Unless otherwise noted, all models, **including baselines**, use GPT-4.1-mini as the summarizer during training and GPT-4.1 during evaluation with top- $k = 10$, while mathematical tasks are summarized by Qwen-3-32B. For evaluation, we sample 500 instances from each benchmark. We use $G = 16$ rollouts. For fair comparison, we train a Qwen-2.5-7B-Instruct summarizer for both training and evaluation on reasoning tasks. Open-domain QA tasks are excluded, as summarization is already an integral component of both our method and the baselines. 7B and 8B models use training batch size 256, rollout batch size 64; 3B doubles both. Learning rate = $2e-6$. Up to 4 retrieval turns are allowed. All models are trained for up to 2 epochs on $8 \times A100$ GPUs (~ 20 hours / 160 GPU-hours). See Appendix C.2 and C.4 for details.

³<https://huggingface.co/BAAI/bge-large-en-v1.5>

⁴<https://nlp.stanford.edu/projects/hotpotqa/enwiki-20171001-pages-meta-current-withlinks-abstracts.tar.bz2>

Models	Types	Methods	Hotpot [†]		2Wiki [†]		Bamb. [‡]		MusiQ. [‡]		Avg	
			F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ
GPT-4.1-mini	Vanilla Methods	CoT	43.7	59.2	48.6	60.8	59.2	76.0	28.3	35.4	45.0	57.9
		Standard RAG	54.5	<u>74.4</u>	41.3	52.4	46.4	51.2	21.9	28.4	41.0	51.6
	Advanced RAG	Self-Ask	65.4	75.0	52.7	57.4	71.7	<u>75.2</u>	31.6	<u>35.0</u>	55.4	60.7
		RAT	56.9	64.2	45.7	49.0	60.3	62.4	29.0	31.4	48.0	51.8
		Search-o1	53.1	74.0	44.4	<u>60.6</u>	<u>63.7</u>	71.2	28.6	33.4	47.5	<u>59.8</u>
Qwen-2.5-7B	Vanilla Methods	CoT	24.9	31.0	25.1	27.6	41.3	43.2	14.8	12.2	26.5	28.5
		Standard RAG	49.2	62.8	32.8	37.6	38.9	40.0	14.4	14.6	33.8	38.8
	RAG-RL	R1-Searcher	<u>71.8</u>	78.0	57.9	63.6	56.5	53.6	33.3	32.6	54.8	57.0
		Search-R1	72.4	78.8	61.0	63.8	58.9	56.8	32.2	32.0	56.1	57.9
		R1-Searcher++	59.0	64.2	<u>61.2</u>	<u>64.4</u>	60.8	59.2	33.8	32.8	53.7	55.2
		ZeroSearch	46.0	50.4	38.4	38.6	35.8	38.4	14.7	13.8	33.7	35.3
	Our Implementations	Vanilla RL	70.9	78.8	<u>61.2</u>	62.4	<u>63.3</u>	63.2	34.4	34.4	<u>57.5</u>	<u>59.6</u>
		UR ²	71.2	79.4	62.6	65.2	64.5	<u>62.4</u>	35.8	34.6	58.5	60.4
Qwen-2.5-3B	Vanilla Methods	CoT	26.6	27.2	22.7	22.6	31.2	33.6	11.3	9.6	23.0	23.3
		Standard RAG	50.6	57.0	29.8	30.4	26.1	27.2	9.7	7.4	29.1	30.5
	RAG-RL	Search-R1	63.1	69.2	49.5	53.4	48.3	48.0	27.6	27.8	47.1	49.6
		Zero-Search	42.7	45.8	26.1	27.6	32.4	31.2	16.9	17.0	29.5	30.4
	Our Implementations	Vanilla RL	65.9	<u>73.6</u>	54.9	58.0	59	<u>57.6</u>	30.0	29.6	<u>52.5</u>	<u>54.7</u>
		UR ²	67.7	76.0	55.2	58.6	57.8	58.4	30.5	31.6	55.3	56.2
LLaMA-3.1-8B	Vanilla Methods	CoT	28.6	31.6	16.4	17.8	43.0	42.4	9.8	10.8	24.5	25.7
		Standard RAG	47.5	54.4	26.2	26.4	26.5	28.0	10.1	10.2	27.6	29.8
	RAG-RL	R1-Searcher	70.8	76.8	59.6	62.2	64.7	62.4	31.1	29.4	56.6	57.7
	Our Implementations	Vanilla RL	70.0	<u>77.6</u>	61.2	64.2	60.6	63.2	<u>32.7</u>	31.8	56.1	<u>59.2</u>
UR ²		<u>70.1</u>	78.8	<u>60.1</u>	<u>63.2</u>	<u>60.7</u>	63.2	34.3	34.0	<u>56.3</u>	59.8	

Table 2: Performance on open-domain QA tasks. We report F1 and LLM-as-a-judge (LSJ) scores, both in %. † = in-domain; ‡ = out-of-domain.

5 Experimental Results

Our UR² framework achieves SoTA performance across reasoning and retrieval tasks, enabling 7B models to match or exceed the GPT model family while significantly outperforming existing RAG and RL-based methods. More comprehensive baseline results can be found in Appendix B.8.

5.1 Main Results on Reasoning Tasks

As shown in Table 1, UR² (GPT Summary) demonstrates substantial improvements across all reasoning tasks on the Qwen-2.5-7B model, achieving average scores of 53.3% on MMLU-Pro, 65.9% on Medicine, and 71.0% on Math benchmarks, representing gains of 3.7%, 5.7%, and 1.2% over the strongest CoT-RL baseline Open-Reasoner-Zero. Using a distilled Qwen-2.5-7B summarizer instead of GPT-4.1 leads to a moderate performance drop, but UR² (qw Summary) consistently outperforms all non-UR² baselines, demonstrating that our gains stem from **the method itself rather than reliance on a strong proprietary summarization model**. Across model scales, UR² shows consistent advantages: on Qwen-2.5-3B,

it achieves even larger performance gains with 9.1% improvement on MMLU-Pro and 8.6% on Medicine over Vanilla RL, demonstrating that UR² provides greater benefits for models with limited knowledge but strong reasoning capabilities. On LLaMA-3.1-8B, it achieves 43.4% on MMLU-Pro, outperforming all baselines. Notably, our method achieves competitive performance with the more capable closed-source GPT-4o-mini model on several tasks. As shown in Tables 1 and 13, advanced RAG methods degrade performance on smaller models and require unacceptable source consumption (except Search-o1).

5.2 Main Results on Open-domain QA

As demonstrated in Table 2, UR² achieves strong performance on open-domain QA, with Qwen-2.5-7B reaching 58.5% average F1 score, outperforming the strongest RAG-RL baseline Search-R1 (56.1%) by 2.4%. UR² demonstrates particularly robust out-of-domain generalization, achieving 64.5% on Bamboogle and 35.8% on MusiQue, surpassing all baselines. Across model scales, UR² maintains consistent advantages: on Qwen-2.5-3B,

Method	MMLU-Pro					Medicine		
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	Avg	MedQA [†]	M-Med [‡]	Avg
UR²	53.2	53.1	72.2	35.0	53.3	69.6	62.8	66.2
w/o Stage-1	48.0	51.1	68.0	30.9	49.5	67.6	63.0	65.3
w/o P_{fallback}	52.0	51.3	68.4	36.6	52.1	71.4	62.0	66.7
w/o Task Mixing	52.2	51.9	68.2	33.2	51.4	70.0	63.6	66.8
w/o LLM Summary	–	–	–	–	–	–	–	–
Vanilla RL	52.2	43.5	64.0	33.8	48.4	64.2	57.4	60.8
4omini Summary	49.3	48.8	67.4	32.4	49.5	65.0	59.2	62.1
Qw3-8B Summary	49.1	49.9	67.8	30.6	49.4	64.8	58.2	61.5
Distilled Qw2.5-7B Summary	49.6	52.3	68.2	31.4	50.4	67.2	61.8	64.5

Table 3: Ablation study of Qwen-2.5-7B-Instruct on MMLU-Pro and medical reasoning tasks. “w/o Task Mixing” means retrieving for all samples. [†] = in-domain; [‡] = out-of-domain.

it achieves 55.3% F1, improving 8.2% over Search-RL; on LLaMA-3.1-8B, it reaches 56.3%, competitive with specialized RAG-RL methods. Notably, our 7B model surpasses GPT-4.1-mini (55.4%) by 3.1%, demonstrating UR²’s effective dynamic coordination of retrieval and reasoning.

6 Analysis

Additional results are provided in the Appendix, including further ablation studies (Appendix B.1), the impact of LLM summaries and corpus design (Appendix B.2), comparative analysis of retrieval integration in RL training (Appendix B.3), analysis of retrieval behavior (Appendix B.4), unsuccessful attempts on reasoning models (Appendix B.5), the compression ratio of the summarization module (Appendix B.6), training efficiency and latency analysis (Appendix C.10), and illustrative case studies (Appendix E).

6.1 Ablation Study

Effect of the Summarization Model. We first examine the role of the summarization model used during training. As shown in Table 3, removing LLM-based summarization (*w/o LLM Summary*) leads to complete training failure, where the model collapses into pure chain-of-thought reasoning and abandons retrieval entirely. This highlights that **retrieval noise is a fundamental challenge in RAG-RL training**, and that effective summarization is crucial for stabilizing retrieval usage. Replacing the default summarizer with alternative or weaker models (4omini, Qw3-8B) results in consistent 3–4% performance drops across tasks, but still substantially outperforms vanilla RL. Notably, the *Distilled Qw2.5-7B Summary* variant exhibits only moderate degradation, demonstrating that UR² does not rely on a strong proprietary summarization model and remains robust under constrained summarization capacity.

Ablation of Training Components. We further analyze ablations on key training components. The *w/o Stage-1* variant shows notable performance drops (e.g., 5.2% in History and 4.2% in Economics), confirming that explicit retrieval activation is essential for learning effective retrieval behavior. Removing the fallback penalty (*w/o P_{fallback}*) slightly improves performance on Law and MedQA, but frequently produces unreasonable or ill-posed queries such as “**which option is right**”, indicating the necessity of regulating fallback behavior. The *w/o Task Mixing* variant yields only minor performance changes, suggesting that selective retrieval primarily improves computational efficiency while maintaining or slightly enhancing accuracy. Overall, these results show that the two-stage optimization and carefully designed reward components jointly contribute to stable and effective retrieval-augmented reasoning.

6.2 Preservation of Intrinsic Reasoning Ability

Method	MedQA	M-Med	Math500	Minerva
CoT	57.2	52.0	76.6	59.4
Vanilla-RL	64.2	57.4	78.2	59.4
UR ² -CoT	59.8	58.8	77.1	61.9
UR ²	69.6	62.8	80.9	61.0

Table 4: Intrinsic reasoning under direct prompts (no retrieval). UR² preserves reasoning ability while retrieval further boosts performance.

A key concern is whether retrieval-augmented training degrades standalone reasoning. We evaluate UR² (Qwen-2.5-7B) under direct reasoning prompts identical to Vanilla-RL, without retrieval. As shown in Table 4, UR²-CoT outperforms CoT and matches Vanilla-RL across knowledge-intensive and reasoning-intensive benchmarks, confirming that difficulty-aware curriculum learning preserves and strengthens intrinsic reasoning.

Corpus	MMLU-Pro					Medicine			Math		
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	Avg	MedQA [†]	M-Med [‡]	Avg	Math500 [†]	Minerva [‡]	Avg
Local Corpus	53.2	53.0	72.2	35.0	53.3	69.6	62.8	65.9	80.9	61.0	71.0
Online Search	57.7	57.8	71.0	35.0	55.4	70.4	65.4	67.9	78.7	61.2	70.0

Corpus	Hotpot [†]		2Wiki [†]		Bamb. [‡]		MusiQ. [‡]		Avg	
	F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ
Local Corpus	71.2	79.4	62.6	65.0	64.5	62.4	35.8	34.6	58.5	60.4
Online Search	62.0	67.6	75.8	81.8	73.7	76.0	34.9	37.8	61.6	65.8

Table 5: Comparison of UR² Qwen-2.5-7B-Instruct using Local Corpus vs. Online Search across different tasks. † = in-domain; ‡ = out-of-domain.

6.3 Impact of Online Search

To test scalability under online retrieval, we compare local corpus with real-time search (Table 5). Online search yields consistent gains on MMLU-Pro (+2.1% average) and medical tasks (+2.0% average), and substantial improvements on 2Wiki (+13.2% F1) and Bamboogle (+9.2% F1), demonstrating strong generalization to scenarios requiring up-to-date or non-Wikipedia knowledge. The only exceptions are math—where Qwen-3-32B’s parametric knowledge already covers the required formulas—and HotpotQA, where rate limits block access to many gold Wikipedia pages. Notably, our setting does not enforce full top-10 coverage, which introduces noise and better reflects real-world retrieval conditions. The substantial gains on out-of-domain benchmarks (Bamboogle, MusiQue) are particularly noteworthy, as they demonstrate that UR²’s learned retrieval strategies transfer effectively to online environments without any re-training. Overall, these results confirm the robustness of UR² in noisy online environments and validate its practical deployment potential.

6.4 Novelty and Concurrent Work

Concurrent large-scale systems such as WebSailor-V2 (Li et al., 2025a), Beyond Ten Turns (Gao et al., 2025), and Tongyi DeepResearch (Tongyi DeepResearch Team, 2025) independently adopt RL-driven multi-step retrieval pipelines with **intermediate summarization**, typically relying on 32B+ models and extensive compute. Their architectural convergence suggests that coordinating retrieval and reasoning under RL is a *structural challenge* rather than an ad-hoc engineering choice. Our contribution is to systematically study and address this challenge, specifically the retrieval–reasoning imbalance, under limited model capacity, showing that effective coordination can be learned even with 3–8B models. As shown in Table 12, UR² learns

task-adaptive retrieval strategies: invoking retrieval sparingly for math (0.94 calls/query) but frequently for knowledge-intensive multi-hop QA (2.46 calls/query), with detailed analysis in Appendix B.4. Unlike these concurrent systems that require extensive infrastructure, UR² achieves comparable or superior performance with a single training run of ~ 160 GPU-hours on $8 \times A100$ GPUs, making the approach accessible to researchers with limited compute budgets.

7 Conclusion

We presented UR², a unified framework that integrates retrieval-augmented generation with reasoning through reinforcement learning. Unlike existing RAG-RL approaches limited to specific domains, UR² demonstrates versatility across mathematical reasoning, medical QA, and open-domain tasks. Our key innovations—difficulty-aware curriculum learning and LLM-summarized retrieval—enable dynamic retrieval-reasoning coordination by learning *when* and *what* to retrieve based on problem difficulty, while preserving native reasoning capabilities. Extensive experiments on Qwen-2.5-3B/7B and LLaMA-3.1-8B show that UR² consistently outperforms advanced RAG and RL baselines, with 7B models achieving performance comparable to GPT-4o-mini and GPT-4.1-mini on several benchmarks. The framework’s robustness across online and offline retrieval settings, diverse corpus configurations, and multiple summarization models confirms its practical applicability. UR² represents a significant step toward adaptive AI systems that flexibly combine parametric knowledge with dynamic information access.

Acknowledgements

This work is supported by the National Natural Science Foundation of China (62372260, 62276152), and Wuxi Research Institute of Applied Technolo-

gies, Tsinghua University. Weizhi Ma is also supported by Beijing Nova Program.

Limitations

While UR² shows strong performance across diverse tasks, some limitations remain. First, our primary experiments use models up to 8B parameters. A preliminary experiment on Qwen-2.5-14B-Instruct using medical training data without summarization (top- $k=2$) shows that UR² still triggers retrieval effectively, achieving 72.4% on MedQA versus 69.8% for CoT, suggesting the framework scales to larger models; however, comprehensive evaluation at larger scales remains for future work. Second, our reliance on on-the-fly LLM-summarized corpora may not fully reflect the complexity of raw web content. Despite these issues, UR² achieves substantial gains (up to 29.6% improvement) and generalizes well across domains.

Future work will explore scaling UR² to 14–32B parameters with full summarization support, incorporating online corpora during training to better capture real-world retrieval dynamics, and investigating more efficient training strategies to reduce costs.

References

- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, and 1 others. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
- Zhipeng Chen, Yingqian Min, Beichen Zhang, Jie Chen, Jinhao Jiang, Daixuan Cheng, Wayne Xin Zhao, Zheng Liu, Xu Miao, Yang Lu, and 1 others. 2025. An empirical study on eliciting and improving rl-like reasoning models. *arXiv preprint arXiv:2503.04548*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*.
- Jiaxuan Gao and 1 others. 2025. Beyond ten turns: Unlocking long-horizon agentic search with large-scale asynchronous rl. *arXiv preprint arXiv:2508.07976*.
- Tianyu Gao, Zhiyuan Liu, Xu Han, Ningyu Zhang, Jilin Tang, Maosong Sun, and 1 others. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2302.05100*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Yiduo Guo, Zhen Guo, Chuanwei Huang, Zi-Ang Wang, Zekai Zhang, Haofei Yu, Huishuai Zhang, and Yikang Shen. 2025b. Synthetic data rl: Task definition is all you need. *arXiv preprint arXiv:2505.17063*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Jian Hu. 2025. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*.
- Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. 2024. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*.
- Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. 2025a. Openreasoner-zero: An open source approach to scaling up reinforcement learning on the base model. *arXiv preprint arXiv:2503.24290*.
- Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. 2025b. GRAG: Graph retrieval-augmented generation. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 4145–4157, Albuquerque, New Mexico. Association for Computational Linguistics.
- Jerry Huang, Siddarth Madala, Risham Sidhu, Cheng Niu, Hao Peng, Julia Hockenmaier, and Tong Zhang. 2025. Rag-rl: Advancing retrieval-augmented generation via rl and curriculum learning. *arXiv preprint arXiv:2503.12759*.

- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*.
- 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*.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Ezgi Korkmaz. 2024. Understanding and diagnosing deep reinforcement learning. *arXiv preprint arXiv:2406.16979*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, and 1 others. 2022. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857.
- Kuan Li and 1 others. 2025a. Websailor-v2: Bridging the chasm to proprietary agents via synthetic data and scalable reinforcement learning. *arXiv preprint arXiv:2509.13305*.
- Weitao Li, Junkai Li, Weizhi Ma, and Yang Liu. 2024. Citation-enhanced generation for LLM-based chatbots. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1451–1466, Bangkok, Thailand. Association for Computational Linguistics.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025b. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*.
- Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhui Chen. 2025. General-reasoner: Advancing llm reasoning across all domains. *arXiv preprint arXiv:2505.14652*.
- Prabhat Nagarajan, Garrett Warnell, and Peter Stone. 2018. Deterministic implementations for reproducibility in deep reinforcement learning. *arXiv preprint arXiv:1809.05676*.
- Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921*.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, and 1 others. 2021. Kilt: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5687–5711.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE.
- Jingyi Ren, Yekun Xu, Xiaolong Wang, Weitao Li, Weizhi Ma, and Yang Liu. 2025. Effective and transparent rag: Adaptive-reward reinforcement learning for decision traceability. *arXiv preprint arXiv:2505.13258*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-math: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025a. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *arXiv preprint arXiv:2503.05592*.
- Huatong Song, Jinhao Jiang, Wenqing Tian, Zhipeng Chen, Yuhuan Wu, Jiahao Zhao, Yingqian Min, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025b. R1-searcher++: Incentivizing the dynamic knowledge acquisition of llms via reinforcement learning. *arXiv preprint arXiv:2505.17005*.
- Hao Sun, Hengyi Cai, Bo Wang, Yingyan Hou, Xiaochi Wei, Shuaiqiang Wang, Yan Zhang, and Dawei Yin.

2024. Towards verifiable text generation with evolving memory and self-reflection. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8211–8227, Miami, Florida, USA. Association for Computational Linguistics.
- Hao Sun, Zile Qiao, Jiayan Guo, Xuanbo Fan, Yingyan Hou, Yong Jiang, Pengjun Xie, Fei Huang, and Yan Zhang. 2025a. Zerosearch: Incentivize the search capability of llms without searching. *arXiv preprint arXiv:2505.04588*.
- Yifan Sun, Jingyan Shen, Yibin Wang, Tianyu Chen, Zhendong Wang, Mingyuan Zhou, and Huan Zhang. 2025b. Improving data efficiency for llm reinforcement fine-tuning through difficulty-targeted online data selection and rollout replay. *arXiv preprint arXiv:2506.05316*.
- Tongyi DeepResearch Team. 2025. Tongyi deepresearch technical report. *arXiv preprint arXiv:2510.24701*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. **MuSiQue: Multi-hop questions via single-hop question composition**. *Transactions of the Association for Computational Linguistics*, 10:539–554.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024a. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290.
- Zihao Wang, Anji Liu, Haowei Lin, Jiaqi Li, Xiaojian Ma, and Yitao Liang. 2024b. Rat: Retrieval augmented thoughts elicit context-aware reasoning in long-horizon generation. *arXiv preprint arXiv:2403.05313*.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256.
- Guangzhi Xiong, Qiao Jin, Xiao Wang, Yin Fang, Haolin Liu, Yifan Yang, Fangyuan Chen, Zhixing Song, Dengyu Wang, Minjia Zhang, and 1 others. 2025. Rag-gym: Optimizing reasoning and search agents with process supervision. *arXiv preprint arXiv:2502.13957*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. *arXiv e-prints*, pages arXiv–2412.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380.
- Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. 2025. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*.
- Xiaoying Zhang, Hao Sun, Yipeng Zhang, Kaituo Feng, Chao Yang, and Helen Meng. 2025. Critique-grpo: Advancing llm reasoning with natural language and numerical feedback. *arXiv preprint arXiv:2506.03106*.
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*.

A Use of AI Tools

In preparing this work, we used commercial LLMs (e.g., Claude 4.0) for non-creative assistance such as language polishing, formatting, and minor code edits. These tools were not involved in method design, experimental setup, or any substantive creative contribution.

B Additional Experimental Results

B.1 Additional Ablation Results

Method	Hotpot [†]	2Wiki [†]	Bamb. [‡]	MusiQ. [‡]
UR ²	71.2	62.6	64.5	35.8
w/o Math Data	70.9	61.2	63.3	34.4
w/o LLM Summary	71.0	62.0	62.7	33.8
weaker Stage-1	69.5	62.2	59.0	34.4

Table 6: Ablation Study of Qwen-2.5-7B-Instruct on open-domain QA. We report F1 scores (in %) here. The second variant removes LLM Summary on top of the first variant without Math Data.

We further conduct ablation studies on UR² in open-domain QA tasks (Table 6). The *w/o Math Data* variant shows minimal impact (0.3-1.4% drops), confirming multi-task training preserves QA performance. Additionally removing LLM Summary causes larger drops on out-of-domain tasks (2.0% on MusiQue) while maintaining in-domain performance, indicating LLM-summarized corpus benefits generalization. The *weaker Stage-1* variant shows the largest degradation on Bamboogle (5.5% drop), highlighting proper retrieval

Method	HotpotQA [†]	2Wiki [†]	Bamboogle [‡]	MusiQue [‡]	Math [†]	Minerva [‡]
Raw Data	72.3	61.9	58.9	36.1	79.0	60.5
Filtered Data	71.0	62.0	62.7	33.8	78.2	59.4

Table 7: Training Data Ablation for Qwen-2.5-7B-Instruct (Vanilla RL on Open-domain QA w/o Summary and Math Tasks). We report F1 scores (%) for open-domain QA.

Method	MMLU-Pro					Medicine		
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	Avg	MedQA [†]	M-Med [‡]	Avg
Qwen-2.5-7B								
Vanilla RL MCQ	47.2	<u>46.1</u>	61.8	33.0	47.0	<u>65.6</u>	<u>60.0</u>	<u>62.8</u>
Vanilla RL	<u>52.2</u>	43.5	<u>64.0</u>	<u>33.8</u>	<u>48.4</u>	64.2	57.4	60.8
UR²	53.2	53.0	72.2	35.0	53.3	69.6	62.8	65.9
Qwen-2.5-3B								
Vanilla RL MCQ	<u>42.3</u>	<u>37.1</u>	<u>57.4</u>	<u>25.0</u>	<u>40.6</u>	50.2	45.0	47.6
Vanilla RL	40.7	34.7	55.0	24.6	38.7	<u>51.8</u>	<u>47.6</u>	<u>49.7</u>
UR²	47.8	49.3	63.9	30.0	47.8	59.8	56.8	58.3

Table 8: Ablation study of Vanilla RL on Qwen-2.5-7B-Instruct and Qwen-2.5-3B-Instruct across multiple-choice reasoning tasks.

initialization is crucial for complex multi-hop reasoning. These results validate our design choices contribute meaningfully across diverse task types.

Table 7 validates the effectiveness of our difficulty-aware data selection strategy. Despite using significantly less training data, the filtered dataset achieves comparable or even superior performance, particularly on out-of-domain tasks (e.g., Bamboogle improves from 58.9% to 62.7%). This indicates that reinforcement learning benefits more from high-quality, difficulty-balanced samples than from large-scale unfiltered data. As a result, UR² enables computationally efficient training while maintaining strong generalization across diverse reasoning tasks.

To ensure a fair comparison, we evaluate Vanilla RL MCQ (MMLU-Pro and Medicine tasks), which trains on mixed multiple-choice tasks similar to UR². As shown in Table 8, Vanilla RL MCQ exhibits task-dependent performance: on Qwen-2.5-7B it improves Medicine performance (62.8% vs. 60.8%) but lowers MMLU-Pro scores (47.0% vs. 48.4%), with the reverse trend on 3B models. Despite these gains, UR² consistently outperforms both Vanilla RL variants across all domains and scales, achieving average improvements on MMLU-Pro of 5.9% for 7B models and 9.1% for 3B models, confirming that its advantage arises from the unified retrieval-reasoning framework rather than task mixing alone.

Overall, the ablations confirm that Stage-1 initialization is crucial for complex reasoning, difficulty-aware filtering yields better performance with fewer samples, and task mixing improves effi-

Summarizer	MMLU-Pro (EM %)				AVG
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	
GPT-4.1	<u>53.2</u>	53.1	72.2	35.0	53.4
Qwen-3-32B	52.5	50.9	<u>72.0</u>	33.6	<u>52.3</u>
Qwen-3-8B	52.5	51.5	71.0	32.8	<u>52.0</u>
GPT-4.1-mini	52.5	51.3	69.4	33.6	51.7
GPT-4o-mini	51.8	49.1	67.4	<u>34.5</u>	50.8
Distilled-Qwen-2.5-7B	49.6	<u>52.3</u>	68.2	31.4	50.4
Qwen-2.5-7B-Instruct	53.4	47.4	67.2	32.0	50.0
w/o Summary	52.1	48.3	68.0	32.2	50.2
Vanilla RL	52.2	43.5	64.0	33.8	48.4

Table 9: Ablation on summarizers in UR² (Qwen-2.5-7B-Instruct) on MMLU-Pro. “w/o Summary” uses top-3 documents without summarizing; “Qwen-3-32B” uses top-16 documents; “Qwen-2.5-7B” (instruct) uses top-5 documents.

ciency without accuracy loss. Importantly, **LLM-summarized retrieval highlights the necessity of addressing retrieval noise in RAG-RL methods**, guiding more stable and generalizable reasoning.

B.2 Impact of LLM Summary and Corpus on UR² Performance

Table 9 examines the robustness of UR² across different LLM summary sources. Remarkably, our framework maintains strong performance regardless of the summarization model quality. While GPT-4.1 achieves the best results (53.4% average), even using smaller open-source models like Qwen-3-8B (52.0%) or budget-friendly APIs like GPT-4o-mini (50.8%) yields substantial improvements over Vanilla RL (48.4%). Most notably, the *w/o Summary* variant still achieves 50.2%—demonstrating that our two-stage training and retrieval-aware prompting mechanisms are inherently robust and

Corpus Summ. Models		Hotpot [†]			2Wiki [†]			Bamb. [‡]			MusiQ. [‡]			Avg			
		F1	LSJ	#R	F1	LSJ	#R	F1	LSJ	#R	F1	LSJ	#R	F1	LSJ	#R	
Abs	✓	ZeroSearch	46.0	50.4	0.66	38.4	38.6	0.73	35.8	38.4	0.54	14.7	13.8	0.62	33.7	35.3	0.64
		Search-R1	72.4	78.8	1.92	61.0	63.8	3.16	58.9	56.8	2.58	32.2	32.0	2.92	56.1	57.9	2.64
		R1-Searcher	71.8	78.0	1.93	57.9	63.6	2.17	56.5	53.6	2.02	33.2	32.6	2.33	54.9	57.0	2.11
		UR ²	71.2	79.4	2.22	62.6	65.0	2.72	64.5	62.4	2.30	35.8	34.6	2.61	58.5	60.4	2.46
Abs	✗	ZeroSearch	44.1	47.0	0.64	32.9	31.8	0.66	32.6	35.2	0.52	14.3	11.8	0.61	31.0	31.5	0.61
		Search-R1	65.8	72.4	2.68	41.8	51.6	3.54	44.8	44.8	2.96	25.1	24.1	3.49	44.4	48.2	3.17
		R1-Searcher	69.7	75.2	2.16	56.6	58	2.45	41.7	40.0	2.38	23.7	22.4	2.84	47.9	48.9	2.46
		UR ²	67.6	73.6	1.98	59.1	59.6	2.53	47.5	47.2	2.10	28.2	25.4	2.43	50.6	51.5	2.26
Full	✓	ZeroSearch	44.3	48.8	0.74	36.5	36.8	0.90	46.3	44.8	0.70	19.3	20.0	0.81	36.6	37.6	0.79
		Search-R1	66.0	67.2	2.01	60.6	65.6	3.12	70.0	71.2	2.06	37.8	39.0	2.69	58.6	60.8	2.47
		R1-Searcher	62.9	68.0	1.97	62.5	66.8	2.15	69.0	65.6	1.86	36.7	37.8	2.24	57.8	59.6	2.06
		UR ²	62.6	68.0	2.11	63.3	67.6	2.73	73.0	74.0	2.13	40.4	42.0	2.55	59.8	62.9	2.38
Full	✗	ZeroSearch	39.2	41.6	0.63	34.0	33.8	0.67	34.1	36.0	0.50	13.4	11.8	0.58	30.2	30.8	0.59
		Search-R1	57.4	60.6	2.75	49.2	51.0	3.50	57.6	55.2	2.82	26.9	26.4	3.40	47.8	48.3	3.12
		R1-Searcher	57.6	61.6	2.24	56.0	59.0	2.37	57.5	57.6	2.07	26.8	26.6	2.63	49.5	51.2	2.33
		UR ²	54.6	60.6	2.03	54.5	55.8	2.51	52.6	49.6	2.06	27.8	26.2	2.38	47.4	48.1	2.25

Table 10: Performance of UR² and baselines on open-domain QA datasets across different corpus configurations. Abs denotes corpora based on Wikipedia abstracts, while Full uses full articles. For each corpus, we use top-10 documents with summaries and top-5 without. #R represents the number of successful retrievals per question.

not dependent on expensive summarization models. This flexibility makes UR² practically deployable across various computational budgets while maintaining its effectiveness, confirming the generalizability of our approach beyond specific model configurations.

Table 10 investigates the impact of different corpus configurations on UR²’s performance across open-domain QA tasks. The results reveal several key insights about corpus design choices. First, using Wikipedia abstracts (Abs, released with HotpotQA) versus full articles (Full) shows task-dependent effects: abstracts perform better on easy questions (HotpotQA), while full articles excel on complex reasoning tasks requiring broader context (2Wiki, Bamboogle, MusiQue). Second, the presence of LLM summarization consistently improves performance across all configurations, with average F1 scores increasing by 6.5-10.8% when summaries are applied. Notably, UR² maintains competitive performance even without summaries (50.6% F1 with Abs, 47.4% with Full), substantially outperforming ZeroSearch’s reliance on synthetic content. The retrieval frequency (#R) analysis shows that UR² strategically balances retrieval calls—using fewer retrievals than Search-R1 while achieving superior performance, demonstrating more efficient knowledge utilization.

Table 11 examines corpus selection for domain-specific tasks, comparing general Wikipedia against specialized MedQA textbooks for medical reasoning. The results demonstrate that domain-

specific corpora provide marginal improvements when summarization is applied (70.2% vs. 69.6% on MedQA), but this advantage diminishes without summaries. More importantly, the performance gap between summarized and non-summarized variants is substantial (8.4% on MedQA with Wikipedia), highlighting that effective summarization is more critical than corpus specialization. This finding suggests that UR²’s LLM-summarized approach can effectively bridge the gap between general and specialized knowledge sources, making it practical for deployment across diverse domains without extensive corpus curation.

Corpus	Summ.	Medicine [†]	M-Med [‡]
Wikipedia	✓	69.6	62.8
Textbooks	✓	70.2	63.8
Wikipedia	✗	61.2	59.2
Textbooks	✗	62.0	58.0

Table 11: Ablation study of UR² on the medical reasoning tasks. We compare different corpus (Wikipedia vs. MedQA Textbooks) and the effect of applying summarization. “w/o Summary” uses top-3 retrieved document.

Collectively, Tables 9, 10, and 11 demonstrate robustness of UR² across three dimensions: corpus configuration, domain specialization, and summarization quality. The framework maintains strong performance whether using abstracts or full articles, general or specialized corpora, and expensive or budget-friendly summarizers. Most remarkably,

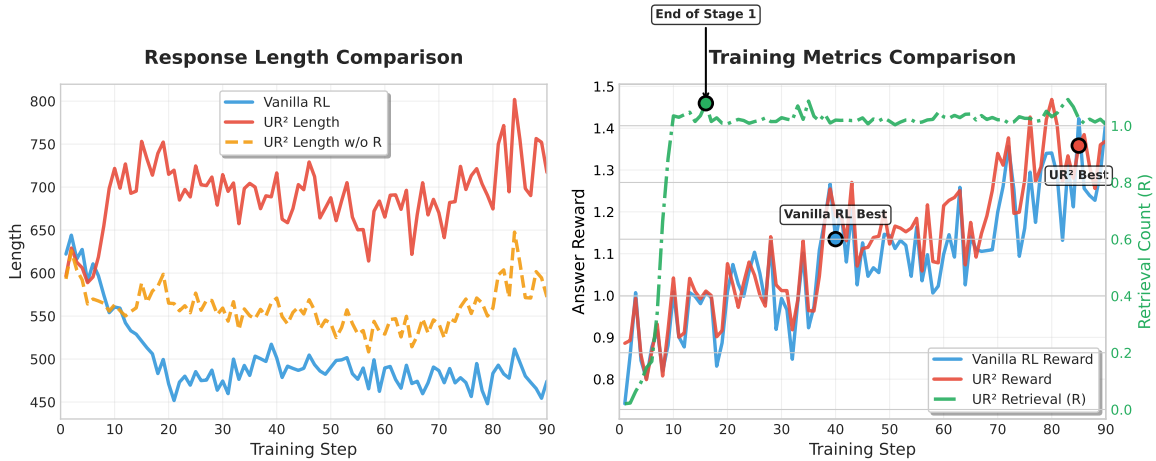


Figure 3: Comparison of Vanilla RL MCQ and UR² performance on Qwen-2.5-3B-Instruct across training steps. Peak test set performances are indicated.

even without any summarization, UR² achieves competitive results through its two-stage training and difficulty-aware retrieval mechanisms. This comprehensive ablation validates that UR²'s effectiveness stems from its fundamental architecture rather than dependency on specific external resources, confirming its practical applicability across diverse computational and domain constraints.

B.3 Comparative Analysis of Retrieval Integration in RL Training

Figure 3 reveals key differences between UR² and Vanilla RL MCQ on Qwen-2.5-3B-Instruct in training. Vanilla RL saturates early at step 40 with 1.1 reward, with later gains mainly due to repeated data every 47 steps. In contrast, UR² steadily improves to 1.4 reward by step 85, matching the 15.7% relative benchmark gain. Retrieval frequency remains dynamic after Stage 1, showing selective use. UR² also generates longer outputs post-training, indicating deeper reasoning. This extended training capability demonstrates that retrieval-augmented approaches fundamentally expand model capacity limits, enabling continuous learning beyond traditional RL saturation points.

B.4 Analysis of Retrieval Behavior

We examine how UR² adapts its retrieval frequency across tasks of varying nature and difficulty. Table 12 reports the average number of retrieval calls per query for UR² Qwen-2.5-7B across different task families.

The retrieval frequency pattern is consistent with task characteristics: mathematical reasoning relies more on internal computation, while multi-hop

Task	Avg. #Retrieval / Query
Math500	0.94
MMLU-Pro	1.34
Open-domain QA	2.46

Table 12: Average retrieval frequency per query for UR² Qwen-2.5-7B across task families. UR² learns to invoke retrieval adaptively based on task characteristics.

QA requires frequent external evidence. Within open-domain QA, retrieval frequency also varies by dataset difficulty: on 2WikiMultiHopQA, which is constructed directly from Wikipedia, retrieval usage remains consistently high across all corpus settings (Table 10), whereas simpler benchmarks trigger fewer calls. Figure 3 further shows that retrieval frequency evolves dynamically during training, stabilizing after Stage 1 as the model learns to selectively invoke retrieval. These results demonstrate that UR² acquires adaptive retrieval strategies rather than relying on fixed heuristics.

B.5 Unsuccessful Attempts on Reasoning Models

We also conducted experiments on the R1-like model DeepSeek-R1-Distill-Qwen-7B⁵. However, when applying the MMLU-Pro prompting setup, we observed that the model lacked any retrieval capability. This remained true even after replacing the original searching special tags with alternative tokens `<search></search>` and `<information></information>`, which were shown in the ablation study (Table 3) to more effectively trigger retrieval. These results indicate a degradation of tool-usage ability after extensive chain-of-thought training. Due to computational

⁵<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B>

constraints, we did not extend training to more updated models such as the Qwen-3 series. We plan to supplement this work with relevant code and experiments in future updates.

B.6 Compression Ratio of the Summarization Module

To quantitatively evaluate the effectiveness of the summarization module, we analyze its compression behavior under both offline and online retrieval settings. Specifically, we randomly sampled 10,000 entries from offline retrieval intermediate files and computed token-level statistics before and after summarization using the Qwen-2.5-7B-Instruct tokenizer.

For offline retrieval, the average number of tokens is reduced from 1,291 to 128 after summarization, corresponding to a compression ratio of $7.87\times$. Although offline retrieval already operates on relatively cleaner and shorter documents, the summarization module still removes a substantial amount of redundant or weakly relevant information, yielding a significantly more compact context for downstream reasoning.

We further conduct the same analysis in a more challenging online retrieval setting, where retrieved documents are typically longer and noisier. Using 3,000 randomly sampled online retrieval instances, we observe that the average token length is reduced from 11,346 to 761, resulting in a much higher compression ratio of $14.91\times$. This stark contrast highlights the growing importance of the summarization module as retrieval environments become more realistic and complex.

Overall, these results demonstrate that the summarization module consistently achieves substantial compression across different retrieval regimes, while its impact is particularly pronounced in online retrieval scenarios. This capability is crucial for controlling context length, reducing computational overhead, and enabling efficient multi-step search-reasoning under real-world conditions.

B.7 Difficulty-Aware Curriculum Preserves Intrinsic Reasoning Ability

We explicitly design the difficulty-aware curriculum to preserve, rather than weaken, the model’s intrinsic reasoning ability. During training, for samples on which the model achieves an accuracy of at least 0.5, we select direct reasoning instead of retrieval-augmented generation. This strategy serves two purposes: (1) it avoids unnecessary re-

trieval calls and reduces API costs, and (2) it ensures that the model continues to practice and retain pure reasoning skills in domains where retrieval is not essential.

As shown in Table 4 (Section 6), UR² consistently outperforms Vanilla-RL on most benchmarks, even when evaluated using direct reasoning prompts. These results demonstrate that difficulty-aware curriculum learning does not degrade intrinsic reasoning ability. Instead, it strengthens the synergy between retrieval usage and standalone reasoning performance.

B.8 Comprehensive Supplementary Results on Open-Domain and Reasoning Tasks

Tables 13 and 14 provide supplementary experimental results, focusing on Advanced RAG methods across different model scales and GPT-4o-mini performance on open-domain QA tasks.

The extended results reveal significant performance limitations of Advanced RAG methods for open-source models. On Qwen-2.5-3B, Self-Ask achieves only 32.0% on MMLU-Pro, substantially underperforming even basic CoT (33.8%). RAT shows inconsistent performance, achieving competitive results on medical tasks (45.0%) but poor performance on Law (18.4%), indicating fragility in cross-domain generalization. Search-o1 demonstrates moderate effectiveness, reaching 41.0% on medical tasks, but fails to achieve consistent improvements across reasoning domains. On LLaMA-3.1-8B, Advanced RAG methods exhibit mixed results. While RAT achieves reasonable performance on Medicine (59.9%) and Math (43.5%), Self-Ask and Search-o1 show notable degradation compared to basic CoT on several sub-domains. These results highlight the challenge of scaling sophisticated retrieval mechanisms to diverse model architectures and reasoning tasks.

GPT-4o-mini establishes strong performance on open-domain QA, with Search-o1 achieving 48.9% F1 average, significantly outperforming other Advanced RAG methods (41.3 and 41.8%). Additionally, RAT and Self-Ask incur prohibitive API costs due to their sentence-level analysis and rewriting operations, making them impractical for large-scale deployment. Notably, Standard RAG achieves competitive performance (42.1% F1) on GPT-4o-mini, suggesting that larger commercial models can effectively leverage simple retrieval without sophisticated coordination mechanisms. The performance gap between GPT-4o-mini (48.9%) and

Method	MMLU-Pro				Medicine			Math			
	Hist. [†]	Phil. [†]	Econ. [†]	Law [‡]	Avg	MedQA [†]	M-Med [‡]	Avg	Math500 [†]	Minerva [‡]	Avg
Qwen-2.5-7B											
<i>Advanced RAG Methods</i>											
Self-Ask	40.7	42.1	60.0	26.2	42.3	51.8	47.8	49.8	74.9	57.7	66.3
RAT	47.2	44.7	64.4	30.0	46.6	60.0	53.2	56.6	74.4	55.5	65.0
Search-o1	42.8	45.9	63.2	29.6	45.4	58.2	52.8	55.6	78.2	60.3	69.3
Qwen-2.5-3B											
CoT	33.6	32.3	48.8	20.6	33.8	39.4	36.8	38.1	63.6	39.9	51.8
Standard RAG	37.8	<u>36.5</u>	51.4	23.2	37.1	45.6	40.0	42.8	65.3	40.8	53.1
<i>Advanced RAG Methods</i>											
Self-Ask	33.1	30.5	44.2	20.2	32.0	39.6	36.0	37.8	58.3	40.3	49.3
RAT	37.8	33.5	50.6	18.4	35.1	47.4	42.6	45.0	67.9	<u>44.9</u>	<u>56.4</u>
Search-o1	33.9	34.5	50.6	20.6	34.9	44.6	37.4	41.0	69.4	43.0	56.2
<i>Our Implementations</i>											
Vanilla RL	<u>40.7</u>	34.7	<u>55.0</u>	<u>24.6</u>	<u>38.7</u>	<u>51.8</u>	<u>47.6</u>	<u>49.7</u>	68.0	43.9	52.4
UR²	47.8	49.3	63.9	30.0	47.8	59.8	56.8	58.3	69.4	45.0	57.2
LLaMA-3.1-8B											
CoT	37.8	40.9	<u>53.4</u>	29.0	<u>40.3</u>	59.6	52.6	56.1	48.4	34.4	41.4
Standard RAG	<u>43.6</u>	33.9	51.0	26.6	38.8	56.4	53.2	54.8	45.0	31.4	38.2
<i>Advanced RAG Methods</i>											
Self-Ask	39.8	32.1	47.0	23.4	35.6	53.0	42.8	47.9	46.9	27.0	37.0
RAT	42.3	37.7	52.6	28.6	40.3	63.8	56.0	59.9	<u>50.1</u>	36.8	43.5
Search-o1	32.6	32.5	46.0	28.0	34.8	56.0	46.0	56.6	41.5	27.8	34.7
<i>Our Implementations</i>											
Vanilla RL	44.6	36.9	53.0	26.4	40.2	<u>66.8</u>	<u>57.4</u>	<u>62.1</u>	45.5	43.4	<u>44.4</u>
UR²	48.3	<u>38.6</u>	58.0	<u>28.8</u>	43.4	68.6	58.4	63.5	54.5	<u>39.0</u>	46.8

Table 13: Extended results on GPT-4o-mini, Qwen-2.5-3B-Instruct, and LLaMA-3.1-8B-Instruct across reasoning tasks. We report EM scores (%) for MMLU-Pro and MedQA, and LLM-as-a-judge scores (%) for math benchmarks. † = in-domain; ‡ = out-of-domain.

smaller models, such as Qwen-2.5-3B (27.8%) for Search-o1, highlights the substantial challenge of achieving effective retrieval-reasoning integration in resource-constrained settings and validates the necessity of our specialized framework design.

C Training Details

C.1 Distilled Qwen-2.5-7B Summarizer Training Setup

To remove reliance on a strong proprietary summarization model, we distill a Qwen-2.5-7B-Instruct summarizer and use it consistently during both training and evaluation.

Training Data. The summarizer is trained on approximately 360 intermediate summaries, including 180k samples over 2 epochs from reasoning tasks (20k per MMLU-Pro subject, 30k from Medicine QA, 30k for Math and 60k for Open-domain QA). All summaries are generated by a mixture of strong teacher models, including GPT-4.1-mini, GPT-4.1, and Qwen-3-32B, to ensure diversity and robustness across domains.

Training Configuration. We perform supervised fine-tuning using a fixed learning rate of 1×10^{-5} . The training runs for 312 GPU hours on H20 GPUs (8×39 hours). No reinforcement

learning is applied during summarizer training.

Usage in UR². Unless otherwise specified, this distilled summarizer is used only for ablation experiments on reasoning tasks, where it replaces proprietary summarizers during both training and evaluation. The corresponding results are reported in the main experiments on reasoning tasks.

C.2 Training Setting Details

We train UR² using the REINFORCE++ algorithm (Hu, 2025), a simplified variant of Proximal Policy Optimization (PPO) designed to encourage exploration. In particular, we discard the critic and omit the KL-divergence term, following previous findings (Zhang et al., 2025; Song et al., 2025a; Chen et al., 2025) that excessive regularization can impede effective strategy learning in sparse-reward scenarios. We retain the PPO-style clipped surrogate objective with $\epsilon = 0.2$ to stabilize policy updates.

The training objective is defined as:

$$J_{\text{UR}^2}(\theta) = \mathbb{E}_{x, \{y^i\}} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|y^i|} \sum_{t=1}^{|y^i|} y_t^i \cdot r_{i,t} \cdot \hat{A}_{i,t} \right] \quad (3)$$

Models	Types	Methods	Hotpot [†]		2Wiki [†]		Bamb. [‡]		MusiQ. [‡]		Avg		
			F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ	F1	LSJ	
GPT-4o-mini	Vanilla Methods	CoT	46.5	51.2	35.0	35.4	55.2	62.4	24.9	26.8	40.4	44.0	
		Standard RAG	<u>59.6</u>	<u>69.6</u>	<u>43.0</u>	<u>45.8</u>	46.7	46.4	19.3	21.6	<u>42.1</u>	<u>45.9</u>	
	Advanced RAG	Self-Ask	45.0	50.4	36.9	40.0	59.3	<u>57.6</u>	<u>26.1</u>	<u>27.8</u>	41.8	44.0	
		RAT	53.8	59.2	34.1	34.8	53.0	51.2	24.3	24.8	41.3	42.5	
		Search-o1	64.3	73.4	47.3	52.0	54.5	56.0	29.6	30.2	48.9	52.9	
Qwen-2.5-7B	Vanilla Methods	CoT	24.9	31.0	25.1	27.6	41.3	43.2	14.8	12.2	26.5	28.5	
		Standard RAG	49.2	62.8	32.8	37.6	38.9	40.0	14.4	14.6	33.8	38.8	
	Advanced RAG	Self-Ask	28.8	61.0	22.2	45.4	28.9	42.4	13.6	19.6	23.4	42.1	
		RAT	37.9	40.6	23.3	23.6	31.6	30.4	14.4	12.4	26.8	26.8	
		Search-o1	50.9	61.6	45.2	48.6	37.5	39.2	20.6	19.8	38.6	42.3	
	RAG-RL	R1-Searcher	<u>71.8</u>	78.0	57.9	63.6	56.5	53.6	33.3	32.6	54.8	57.0	
		Search-R1	72.4	78.8	61.0	63.8	58.9	56.8	32.2	32.0	56.1	57.9	
		R1-Searcher++	59.0	64.2	<u>61.2</u>	<u>64.4</u>	60.8	59.2	33.8	32.8	53.7	55.2	
		ZeroSearch	46.0	50.4	38.4	38.6	35.8	38.4	14.7	13.8	33.7	35.3	
	Our Implementations	Vanilla RL	70.9	78.8	<u>61.2</u>	62.4	63.3	63.2	34.4	34.4	57.5	59.6	
		UR ²	71.2	79.4	62.6	65.2	64.5	<u>62.4</u>	35.8	34.6	58.5	60.4	
Qwen-2.5-3B	Vanilla Methods	CoT	26.6	27.2	22.7	22.6	31.2	33.6	11.3	9.6	23.0	23.3	
		Standard RAG	50.6	57.0	29.8	30.4	26.1	27.2	9.7	7.4	29.1	30.5	
	Advanced RAG	Self-Ask	33.8	47.2	21.0	28.8	30.6	32.0	14.5	14.8	25.0	30.7	
		RAT	30.1	32.2	15.1	15.4	30.6	28.0	11.0	8.2	21.7	21.0	
		Search-o1	36.4	37.6	30.8	31.8	31.4	32.0	12.5	10.0	27.8	27.9	
	RAG-RL	Search-R1	63.1	69.2	49.5	53.4	48.3	48.0	27.6	27.8	49.6	47.1	
		Zero-Search	42.7	45.8	26.1	27.6	32.4	31.2	16.9	17.0	29.5	30.4	
	Our Implementations	Vanilla RL	<u>65.9</u>	<u>73.6</u>	<u>54.9</u>	<u>58.0</u>	<u>59</u>	<u>57.6</u>	<u>30.0</u>	<u>29.6</u>	<u>52.5</u>	<u>54.7</u>	
		UR ²	67.7	76.0	55.2	58.6	57.8	58.4	30.5	31.6	55.3	56.2	
	LLaMA-3.1-8B	Vanilla Methods	CoT	28.6	31.6	16.4	17.8	43.0	42.4	9.8	10.8	24.5	25.7
			Standard RAG	47.5	54.4	26.2	26.4	26.5	28.0	10.1	10.2	27.6	29.8
Advanced RAG		Self-Ask	43.0	50.8	27.3	29.8	41.5	44.8	16.8	16.4	32.2	35.5	
		RAT	44.5	48.8	16.4	15.6	39.7	39.2	17.0	16.0	29.4	29.9	
		Search-o1	53.0	59.4	37.5	38.4	30.0	30.4	15.9	16.2	34.1	36.1	
RAG-RL		R1-Searcher	70.8	76.8	59.6	62.2	64.7	62.4	31.1	29.4	56.6	57.7	
Our Implementations		Vanilla RL	70.0	<u>77.6</u>	61.2	64.2	60.6	63.2	<u>32.7</u>	<u>31.8</u>	56.1	<u>59.2</u>	
		UR ²	<u>70.1</u>	78.8	<u>60.1</u>	<u>63.2</u>	<u>60.7</u>	63.2	34.3	34.0	<u>56.3</u>	59.8	

Table 14: Extended results of GPT-4o-mini, Qwen-2.5-3B, and LLaMA-3.1-8B on open-domain QA. We report F1 and LLM-as-a-judge (LSJ) scores, both in %. † denotes in-domain datasets; ‡ indicates out-of-domain.

where the importance weight is:

$$r_{i,t} = \frac{\pi_{\theta}(y_t^i | x, y_{<t}^i; o_i)}{\pi_{\text{old}}(y_t^i | x, y_{<t}^i; o_i)} \quad (4)$$

and the normalized advantage is:

$$\hat{A}_{i,t} = \text{Norm}_{\text{batch}}(\text{Norm}_{\text{group}}(R_i - b)) \quad (5)$$

The advantage $\hat{A}_{i,t}$ is computed by subtracting the group-level reward baseline and applying normalization across the group and batch to improve learning stability. Here, x denotes the input prompt, $\{y^i\}$ are the sampled trajectories, o_i is the retrieved context and b is the group-level baseline (mean of R_i).

To reduce overfitting to retrieved content, we adopt a retrieval masking strategy (Sun et al.,

2025a; Song et al., 2025a; Jin et al., 2025), which treats retrieved external knowledge as part of the observation space rather than trainable input. This encourages the model to reason based on retrieved information without directly optimizing on it. Our implementation builds upon the REINFORCE++ baseline provided by OpenRLHF (Hu et al., 2024).

Each prompt is rolled out $G = 16$ times. We use the mean reward of each rollout group as the baseline for computing the advantage of each sample. To stabilize training, we apply a two-stage normalization scheme: normalization is first performed within each rollout group, followed by global normalization across the full batch.

Training is conducted with DeepSpeed ZeRO-2 (Rajbhandari et al., 2020) for memory efficiency.

We use gpt-4.1-mini-2025-04-14 as the summarization model during training. Token limits per generation turn are set to 3072 for math tasks, 1536 for multiple-choice questions (MCQ), and 512 for open-domain QA. Sampling parameters are fixed as temperature = 1.0 and top_p = 0.9.

We train for up to 2 epochs. In practice, most models achieve optimal performance within 1.5 epochs. Therefore, we report results from the checkpoint with the best test set performance within the first 1.5 epochs. We save checkpoints every 5 steps for single-task training and 3B models, and every 10 steps for larger-scale experiments. The specific training steps for each reported model are detailed in Table 15 below. *W/o Stage-1* variant in Table 3 replaces the special tags with `<search></search>` and `<information></information>`, removing the initial retrieval capability activation stage. The *Weaker Stage-1* variant in Table 6 employs a modified training protocol based on UR² Qwen-2.5-7B-Instruct for MCQ tasks, where retrieval-related rewards are only provided during the initial 10 training steps. The *Qw3-8B* variant in Table 3 uses Qwen-3-8B for summarization with max_tokens = 2048, temperature = 0.3, and top_p = 0.7. Specifically, the retrieval reward assigns 0.5 for single retrieval attempts and 1.0 for multiple retrievals (≥ 2), reflecting a more conservative retrieval activation strategy than that of our proposed method.

Table 15 summarizes the training configurations and checkpoint details across all model scales. Two key observations can be drawn:

For the task mixing strategy of multiple-choice reasoning tasks, we combine MedQA and synthetic MMLU-style data, assigning most hard questions to retrieval-augmented training while maintaining an overall 1:1 ratio between retrieval-based and direct reasoning instances.

First, for Qwen models, performance consistently improves as more training compute is introduced via our UR² method. The method’s design—encouraging structured retrieval behavior—ensures that increased steps and effective epochs lead to meaningful gains across tasks.

Second, while UR² also improves performance on LLaMA-3.1-8B, training on this model is observed to be less stable. Performance tends to saturate early (e.g., low effective epochs despite higher step counts), for both Vanilla RL and UR² variants. This indicates that LLaMA-3.1-8B may require different training strategies to maintain learning dy-

namics over time. Future work will explore alternative foundation models and optimization schedules to improve convergence and stability.

C.3 Vanilla RL

To ensure a fully fair comparison, Vanilla RL is a baseline following the same setup and datasets as UR², applying RAG-RL to open-domain QA and CoT-RL to math/MCQ tasks. Specifically: (1) identical or subset training data, (2) REINFORCE++ algorithm, (3) same hyperparameters (KL coef, batch size), and (4) standard domain-specific training methods (CoT-RL in reasoning tasks and RAG-RL in open-domain QA). In the retrieval scenario, it serves as a single stage multi-step search–reasoning training baseline similar to Search-R1, using the same setup and subset training data as UR². In the reasoning scenario, it serves as a single stage onestep reasoning training baseline similar to SimpleRL-Zoo, using the same setup and identical or subset training data as UR².

C.4 Evaluation Details

All evaluations are performed using vLLM version 0.6.5. The vLLM version of Qwen-3 used is 0.8.5.post1. In evaluation, We maintain the same max_tokens limits used during training: 3072 for math benchmarks, 1536 for MCQ, and 512 for open-domain QA per generation step. For GPT-family models, these limits are increased to 4096, 2048, and 1024, respectively. For sampling during evaluation, we use more conservative hyperparameters: temperature = 0.3 and top_p = 0.5, aiming for higher answer consistency. Summarization for math tasks is conducted using Qwen-3-32B with max_tokens = 8192, temperature = 0.3, and top_p = 0.7. Final evaluation summarization is performed using gpt-4.1-2025-04-14 with max_tokens = 2048, temperature = 0.3, and top_p = 0.5.

The RL methods mentioned in this paper all follow the settings described in their original works. Specifically, Open-Reasoner-Zero, General Reasoner, SimpleRL-Zoo, R1-Searcher, Search-R1, and ZeroSearch are implemented using the Qwen-2.5-Base models. Although an Instruct version of Search-R1 exists, its performance is significantly inferior and thus excluded from comparison. R1-Searcher with LLaMA-3.1-8B adopts the Instruct variant. Vanilla methods, including CoT and standard RAG, are applied using the Instruct versions for all open-source models.

Model	Training Dataset	Dataset Size	Checkpoint Step	Training Epochs
Qwen-2.5-3B - Main Experiments				
UR ² -Math&QA	Math&QA	6000	47	1.0
UR ² -MCQ	MMLU&Medqa	9000	85	1.2
Vanilla RL-Math	Math	3000	15	0.64
Vanilla RL-QA	QA	3000	40	1.7
Vanilla RL-MMLU	MMLU	6000	47	1.0
Vanilla RL-MedQA	Medqa	3000	40	1.7
Vanilla RL-MCQ	MMLU&Medqa	9000	40	0.57
LLaMA-3.1-8B Models - Main Experiments				
UR ² -Math&QA	Math&QA	6000	30	0.32
UR ² -MCQ	MMLU&Medqa	9000	30	0.21
Vanilla RL-Math	Math	3000	30	0.64
Vanilla RL-QA	QA	3000	47	1.0
Vanilla RL-MMLU	MMLU	6000	60	0.64
Vanilla RL-MedQA	Medqa	3000	30	0.64
Qwen-2.5-7B - Main Experiments				
UR ² -Math&QA	Math&QA	6000	40	0.43
UR ² -MCQ	MMLU&Medqa	9000	100	0.71
Vanilla RL-Math	Math	3000	40	0.43
Vanilla RL-QA	QA	3000	25	0.53
Vanilla RL-MMLU	MMLU	6000	94	1.0
Vanilla RL-MedQA	Medqa	3000	47	1.0
Vanilla RL-MCQ	MMLU&Medqa	9000	60	0.43
7B Models - Ablation Studies				
Ablation-MCQ-w/o P_{fallback}	MMLU&Medqa	9000	110	0.78
Ablation-MCQ-w/o Stage-1	MMLU&Medqa	9000	110	0.78
Ablation-MCQ-w/o Task Mixing	MMLU&Medqa	9000	120	0.85
Ablation-MCQ-QW3 summary	MMLU&Medqa	9000	50	0.36
Ablation-MCQ-4omini summary	MMLU&Medqa	9000	50	0.36
Ablation-Math&QA weaker Stage-1	Math&qa	6000	80	0.85
Ablation-QA w/o LLM summary	QA	3000	60	1.28
Ablation-QA Raw data	R1-Searcher	8148	70	0.55
Ablation-Math Raw data	SimpleRL-Zoo	16662	10	0.056

Table 15: Training checkpoint details for UR² models. Checkpoints were saved every 5 steps for 3B models and single-task training, and every 10 steps for larger models. Main experiments use the full training configuration, while ablation studies vary specific components.

Advanced RAG Baseline Implementations:

Search-o1 with Retrieval-Augmented Generation: We adapt the Search-o1 framework (Li et al., 2025b) to operate within a controlled evaluation environment. While maintaining its core iterative reasoning mechanism and document analysis capabilities, our implementation leverages the KILT Wikipedia corpus with BGE-large-en-v1.5 embeddings for knowledge retrieval. This approach consolidates the multi-agent architecture into a unified model with structured prompting, ensuring consistent evaluation across all baselines while preserving the essential reasoning patterns.

Self-Ask with Retrieval-Augmented Generation: Our implementation follows the Self-Ask framework’s (Press et al., 2023) question decomposition strategy, employing batch retrieval from the local KILT corpus to enhance efficiency. The system maintains the characteristic “Follow up:” and “Intermediate answer:” reasoning chain format, with stopping criteria incorporating both semantic completion detection and a maximum of 10 follow-up questions. When decomposition challenges arise, the framework seamlessly transitions to standard RAG, ensuring robust performance across diverse question types.

RAT (Retrieval-Augmented Thought): We adapt RAT (Wang et al., 2024b) for unified evaluation across reasoning and QA tasks. The framework retains the core principle of knowledge-enhanced reasoning while operating at the paragraph level rather than the sentence level, with corresponding modifications to the prompting strategy. This design choice maintains consistency with our evaluation infrastructure while capturing RAT’s fundamental insight of augmenting reasoning processes with relevant external knowledge.

All advanced RAG methods operate within a standardized retrieval infrastructure: documents are retrieved from the 100-word segmented KILT Wikipedia corpus (29M documents in total). For GPT-family models, we use top- $k=10$ retrieval. Due to model limitations, LLaMA and Qwen variants use top- $k=5$. For summarization or other auxiliary operations beyond reasoning, each model performs the processing itself rather than relying on GPT-4.1, ensuring consistency with its own capabilities.

Online Corpus Retrieval Implementation:

To evaluate the generalization capability of UR² with real-world web content, we implement an online corpus retrieval system that dynamically fetches and processes web documents. Unlike the offline Wikipedia corpus used during training, this online retrieval mechanism provides access to up-to-date information from the internet.

The online retrieval pipeline consists of three main components:

Web Search and Content Extraction: We utilize the **Bing Search API** to retrieve relevant URLs based on the model’s search queries. To ensure robust retrieval quality, we implement a multi-round crawling strategy with up to three rounds of attempts. In each round, the system fetches $k \times 3$ candidate URLs and crawls them in parallel using a thread pool with 256 workers. The system implements intelligent retry logic—if the initial k URLs fail to provide sufficient valid content, it automatically attempts to crawl additional URLs from the candidate pool. This approach significantly improves the success rate of obtaining high-quality content.

HTML-to-Markdown Conversion: Raw HTML content from web pages often contains noise such as navigation elements, advertisements, and scripts. We deploy a dedicated service using

ReaderLM-v2-1.5B Model ⁶ through the vLLM framework to convert HTML to clean Markdown format. The preprocessing pipeline removes script tags, style elements, base64-encoded images, and other irrelevant content using optimized regular expressions. The model then generates readable Markdown that preserves the main textual information while discarding formatting artifacts. To improve efficiency, we implement an LRU cache with a capacity of 10,000 entries, achieving significant speedup for repeated content.

Content Summarization: The summarization prompt is carefully designed to distinguish between knowledge-based queries (which can be answered with factual information) and reasoning-based queries (which require complex computation). For knowledge-based queries, the model extracts and presents relevant facts; for reasoning-based queries, it returns a fallback message indicating that direct reasoning is more appropriate. The summarizer here is GPT-4.1-2025-04-14.

The entire pipeline is orchestrated through a FastAPI service that handles concurrent requests efficiently. Rate limiting is enforced for the Bing API (95 requests per second) to comply with usage policies. The system maintains detailed logging for debugging and performance monitoring, tracking metrics such as cache hit rates, crawling success rates, and end-to-end latency.

Due to network and hardware limitations, a small portion of Wikipedia pages failed to be crawled correctly, and a subset of queries did not receive valid responses. Given constraints on time and budget, no additional remediation was applied to these cases. However, this reflects the system’s alignment with real-world deployment settings, where large-scale QA systems must be robust to occasional retrieval failures and operate under imperfect infrastructure.

This online retrieval implementation enables UR² to access current information beyond its training data, demonstrating its ability to integrate real-time knowledge into the reasoning process.

C.5 Training Dataset Details

We construct a unified training set that spans multiple task domains to ensure comprehensive coverage of diverse reasoning and knowledge-based challenges. For mathematical reasoning capabilities, we incorporate data from the training split

⁶<https://huggingface.co/jinaai/reader-lm-1.5b>

of SimpleZoo-RL, which provides a rich collection of mathematical problem-solving scenarios from (Hendrycks et al., 2021; Cobbe et al., 2021). Note that since the original SimpleZoo-RL data is relatively simple, medium- and hard-difficulty questions are largely missing, resulting in an overall easy:medium:hard ratio of 1:1:1 rather than the 7:2:1 used in Section 3.3. Moreover, due to limitations of LLaMA-3.1-8B-Instruct, we substitute easy-difficulty questions for hard ones during training. To enhance open-domain QA performance, we include samples from the R1-Searcher dataset, which spans a broad range of questions derived from the training sets of 2Wiki and HotpotQA. For specialized domain knowledge, particularly in the medical field, we utilize multi-choice questions from MedQA, ensuring our model can handle domain-specific reasoning in healthcare contexts.

To further diversify our training data and extend coverage to humanities subjects, we generate synthetic questions in three additional domains: philosophy, history, and economics. These synthetic questions are created using Qwen-3-32B and follow the MMLU-Pro format to maintain consistency with established academic evaluation standards. Specifically, we use 5-shot prompting with MMLU-Pro development set examples to generate 10 questions with 4–10 options each. We discard format-non-compliant questions and observe the model’s tendency to generate simple questions with 4–5 options, so we request the model to produce additional options and increase the difficulty for each question. For quality control, we use GPT-4o-mini-2024-07-18 to evaluate each question’s correctness three times, discarding any question identified as incorrect in any evaluation. We then employ Qwen-2.5-7B-Instruct for difficulty assessment, finding approximately 80% of questions are easy-level. We randomly sample difficult questions as seeds for subsequent generations, using different seeds for each batch. Given that downstream test sets contain subject subdivisions (e.g., Economics encompasses microeconomics, macroeconomics, and econometrics), we utilize Qwen-3-32B to classify questions by subdomain, ensuring comprehensive coverage. We repeat this pipeline for 3–4 iterations to obtain the final training set.

Notably, our synthetic questions differ from MMLU-Pro in emphasizing multi-hop reasoning rather than specific knowledge points. This is evident of our results in Table 1 where Vanilla RL shows limited improvement over CoT Base-

line for Qwen-2.5-7B-Instruct and LLaMA-3.1-8B-Instruct (3.9% and -0.1% respectively), **demonstrating no overfitting to the test set**. Despite these characteristics, UR² consistently achieves improvements across models, validating our method’s effectiveness.

C.6 About Fallback Fault in Retrieval Corpus Construction

When the policy model generates an invalid search query that triggers a fallback message from the LLM summarizer (i.e., *This query requires design, computation, or complex reasoning, which exceeds the capabilities of a search engine. Please input another query or proceed with direct reasoning.*), we observe that due to the use of retrieval masking, the model gradually learns to treat the content within `<info>...</info>` as informative for reasoning. As a result, when a fallback fault is encountered, the model tends to hallucinate. Therefore, we append the following visible message after `</info>` during training to mitigate this issue: *It seems that this query exceeds the capabilities of the retrieval system. We may consider rephrasing it into a more fact-based and searchable question that does not require complex reasoning, or proceed with direct reasoning based on prior knowledge.*

C.7 Stage 1 Training Details

Due to the involvement of multiple models and tasks, Section 3.4 only presents the stage-1 setup for Qwen-2.5-7B-Instruct on math and open-domain QA. Here, we elaborate on the initialization strategies for other models and tasks.

Math and Open-Domain QA. We use the discarded math training samples with rollout accuracy below 0.2 as cold-start data. These harder examples naturally increase the likelihood of triggering retrieval. For Qwen-2.5-3B-Instruct, its limited capacity makes it more prone to Format violations when invoking retrieval. Since each violation incurs a -1 penalty, the original retrieval reward (+3 for one query, +4 for two or more) becomes insufficient to incentivize retrieval. To address this, we increase the retrieval rewards to +5 and +7, respectively. In contrast, LLaMA-3.1-8B-Instruct tends to retrieve for almost every question in early steps. To prevent over-reliance on retrieval and preserve reasoning ability, we remove the extra reward for multiple queries and assign a fixed +3 reward upon any retrieval activation.

MMLU-Pro and Medicine Tasks. Unlike math tasks, MMLU-Pro and medicine tasks often require domain-specific knowledge, and retrieval is less likely to lead to fallback faults. For LLaMA-3.1-8B-Instruct and Qwen-2.5-7B-Instruct, a weak reward signal is sufficient during early training: +0.5 for one valid retrieval and +1 for two or more. Unlike the original stage-1 design for math and open-domain QA, this version also incorporates answer rewards from the beginning, facilitating early alignment with task-specific correctness (i.e., no longer relying on cold-start data). In this variant, retrieval rewards are only applied during the first 10 training steps and then disabled.

For Qwen-2.5-7B-Instruct trained on math and open-domain QA, we adopt the stage-1 setup originally used for the MMLU-Pro and medicine tasks, corresponding to the weaker Stage-1 variant in Table 6.

For Qwen-2.5-3B-Instruct, we extend Stage 1 to 15 steps. To encourage retrieval, outputs that do not invoke any retrieval call are penalized with a -1 Format Reward (non-accumulative).

C.8 On Randomness and Reproducibility

RL training is known to exhibit inherent instability and variability across runs, often leading to divergent results even under identical settings (Nagarajan et al., 2018; Korkmaz, 2024). This randomness is attributed to factors such as stochastic policy updates, environment interactions, and non-deterministic hardware behavior. Despite these challenges, our experiments demonstrate remarkable stability. Thanks to the incorporation of Batch Normalization and Group Normalization in reward calculation, all models converge successfully in a single training run.

During evaluation and result aggregation, we employed a non-zero temperature setting to maintain controlled output diversity, thereby enhancing performance and mitigating the risk of repetitive generations. Due to the substantial API costs associated with GPT-4.1, conducting multiple evaluation runs to average results was not feasible. Nevertheless, given that the datasets contain approximately 500 samples—providing sufficient statistical power—we performed a targeted reproducibility assessment on HotpotQA using the UR² Qwen 7B-Instruct model. Specifically, three independent evaluation runs yielded F1 scores of 71.7, 71.9, and 71.2, respectively. These consistent results indicate that stochasticity exerts minimal influence on evalua-

tion metrics and comparative model assessments. Furthermore, we conducted supplementary evaluations on all identified outlier cases across baseline and proposed methods to ensure the robustness of our findings.

C.9 API Consumption

We measured the API usage cost of UR² Qwen-2.5-7B-Instruct on MCQ tasks and its w/o Stage-1 variant. Training 100 steps with UR² using GPT-4.1-mini cost approximately \$320, while the w/o Stage-1 variant cost around \$100. Additionally, summarization and testing on HotpotQA using GPT-4.1 for UR² Qwen-2.5-7B-Instruct cost about \$20 per run. Since the training is a one-time expense, we consider the overall training-related consumption acceptable. Furthermore, experiments reported in Section B.2 and Table 3 show that substantial performance gains can be achieved without relying on closed-source models, suggesting that open-source models or less expensive APIs provide a viable alternative for achieving comparable improvements.

C.10 Training Efficiency and Latency

For models up to 8B parameters with a maximum generation length of 2,048 tokens, UR² training requires approximately 20 hours of wall-clock time on $8 \times A100$ GPUs, corresponding to ~ 160 GPU-hours in total. All RL baselines are trained under identical hardware and comparable step budgets, so the performance gains are achieved at no additional deployment cost and similar training budgets.

Regarding inference latency, the summarization module introduces an extra generation step per retrieval call. However, by compressing retrieved context from an average of 11,346 tokens to 761 tokens ($14.91 \times$), it significantly reduces the reasoning model’s input length. Given the quadratic attention cost $O(n^2)$, this compression substantially decreases the downstream reasoning time. In practice, when the summarizer is co-located on the same hardware or served via a fast API, the end-to-end latency of UR² is comparable to—and in long-context scenarios lower than—standard RAG pipelines that feed raw documents directly. Since UR² does not increase model parameters, deployment overhead is identical to the base model.

D Prompts used in Experiments

D.1 Prompts of LLM-as-a-Judge

Prompt for Math Evaluation

Instruction:

You are an expert math evaluator. Given a question, a gold answer and a predicted answer, judge if they are mathematically consistent.

Ignore formatting (e.g., `\text{}`, spacing, capitalization). Accept equivalent expressions (e.g., factored vs expanded form). If the prediction matches only part of a multi-part answer (e.g., one of several intervals or roots), label it as

Partially correct.

Output format:

- Reason: Brief explanation
- Judgment: Correct / Partially correct / Incorrect

Input:

- Question: {question}
 - Gold: {gold}
 - Pred: {pred}
-

Prompt for RAG Evaluation

Instruction:

Given a Question and its Golden Answer, verify whether the Predicted Answer is correct. The prediction is correct if it fully aligns with the meaning and key information of the Golden Answer. Respond with True if the prediction is correct and False otherwise.

Input:

- Question: {question}
- Golden Answer: {gold_answer}
- Predicted Answer: {predicted_answer}

Your response should be exactly "True" or "False"

D.2 Prompts of Evaluation and Training

Prompt for MMLU-Pro&MedQA

Instruction:

You are solving a multiple-choice question. Analyze each option carefully and logically. Think step by step: consider the meaning and implications of each option, eliminate incorrect ones with clear reasoning, and select the best answer through comparison.

During your reasoning, if you're unsure about any fact, you may issue a **search query** like this: `<|begin_of_query|>` your concise query (less than 20 words) `<|end_of_query|>`

- You can issue **multiple queries** at different steps in your reasoning.
- **Each query must target only one fact or statement.** Do not combine multiple ideas in a single query.
- **Examples:**
 - ✓ `<|begin_of_query|>` What are the common symptoms of pneumonia? `<|end_of_query|>`
 - ✓ `<|begin_of_query|>` What is the typical treatment for pneumonia in elderly patients? `<|end_of_query|>`
 - ✗ `<|begin_of_query|>` What are the symptoms and treatments for pneumonia in elderly patients? `<|end_of_query|>`
- You may issue **at most four queries** in total — use them wisely.

Once documents are returned in this format:

`<|begin_of_documents|>` ... (search results here) `<|end_of_documents|>`

Use the retrieved documents to verify, reject, or revise your prior reasoning about the options. Then continue analyzing the options until you're confident in your answer.

Final answer format: the correct answer is: A, B, C, D, etc. (only the letter corresponding to the correct option)

Prompt for Math

Instruction:

You are solving a math problem. Think step by step to solve it.

The reasoning process includes detailed considerations such as analyzing questions, summarizing relevant findings, brainstorming new ideas, verifying the accuracy of current steps, refining any errors, and revisiting previous steps.

During your reasoning, if you're unsure about a factual concept — such as a definition, formula, theorem, or mathematical constant — you may issue a **search query** to clarify it.

Format your query using the following template (each query must target only one fact):

`<|begin_of_query|>` your concise query (less than 20 words) `<|end_of_query|>`

✓ **Examples:**

- `<|begin_of_query|>` Definition of Möbius function `<|end_of_query|>`
- `<|begin_of_query|>` Formula for variance of Bernoulli distribution `<|end_of_query|>`

✗ **Do NOT query for reasoning-related content like:**

- Whether a solution approach is valid
- How to compute a specific value
- Multi-step deductions or conclusions

You may issue at most **four** search queries per problem — use them wisely.

When documents are returned in this format:

`<|begin_of_documents|>` ... (search results here) `<|end_of_documents|>`

Use the evidence to confirm or revise your reasoning. Then continue analyzing the question until you're confident in the answer.

At the end of your reasoning, give your final answer in the following format:

`\boxed{YOUR_ANSWER}`

Prompt for Open-Domain QA

Instruction:

You are solving a factual open-domain question from a Knowledge Question Answering (KQA) task. The question requires step-by-step reasoning over real-world knowledge to identify a specific, factually correct answer.

Carefully analyze the question to understand the key entities, relationships, and constraints involved. Retrieve and consider relevant factual knowledge, and reason logically to identify the most accurate answer.

During your reasoning, if you're unsure about any fact, you may issue a **search query** like this: `<|begin_of_query|>` your concise query (less than 20 words) `<|end_of_query|>`

- You can issue **multiple queries** at different steps in your reasoning.
- **Each query must target only one fact or statement.** Do not combine multiple ideas in a single query.
 - ✓ **Example:**
 - * `<|begin_of_query|>` When did Einstein move to the United States? `<|end_of_query|>`
 - * `<|begin_of_query|>` Why did Einstein leave Germany?
`<|end_of_query|>`
 - ✗ **Do not combine them like this:**
 - * `<|begin_of_query|>` When did Einstein move to the US and why did he leave Germany?
`<|end_of_query|>`
- You may issue **at most five queries** in total — use them wisely.

Once documents are returned in this format:

`<|begin_of_documents|>` ... (search results here) `<|end_of_documents|>`

Use the evidence to confirm or revise your reasoning. Then continue analyzing the question until you're confident in the answer.

At the end of your reasoning, give your final answer in the following format: `\boxed{YOUR_ANSWER}`

D.3 Prompts for Summarizing

Prompt for Summarizing Math Documents During Evaluation

Task Instruction:

You are assisting in solving a math problem. You are tasked with reading and analyzing Wikipedia content based on the following inputs: **Previous Reasoning Steps**, **Current Search Query**, and **Wikipedia Content**. Your task is to extract accurate and relevant information from the provided Wikipedia content to support or enhance the reasoning process.

- Carefully read the provided **Wikipedia Content**;
- Extract factual information that can:
 - Directly assist in answering the **Current Search Query**, or
 - Help validate, complete, or correct earlier reasoning steps.
- The extracted information should be:
 - Accurate and trustworthy;
 - Closely relevant to the query;
 - Helpful in improving, expanding, or supporting the mathematical reasoning.

Important: Do NOT attempt to correct or rewrite the previous reasoning. Treat it only as contextual reference that may be flawed.

Output Format:

Present the information beginning with the label ****Final Information**** as shown below.

****Final Information****

[Helpful factual information]

Inputs:

- Previous Reasoning Steps: {prev_reasoning}
 - Current Search Query: {search_query}
 - Wikipedia Content: {wikipedia_content}
-

Prompt for Summarizing Math Documents During Training

Task Instruction:

You are assisting in solving a math problem. Your task is to determine whether the current query requires external factual knowledge (such as definitions, formulas, theorems, or lookup values), and if so, extract accurate and relevant information from the provided Wikipedia content to support or enhance the reasoning process.

Step 1: Classify the Query Type

Determine whether the query falls into one of the following categories:

- **Knowledge-based query:** Can be directly answered using factual knowledge.
- **Reasoning-based query:** Requires multi-step deduction, logical reasoning, or constructive computation.

If reasoning-based, return: *This query requires design, computation, or complex reasoning, which exceeds the capabilities of a search engine. Please input another query or proceed with direct reasoning.*

Step 2: Analyze Knowledge-Based Queries (if applicable)

- Carefully read the Wikipedia Content;
- Extract factual information that:
 - Directly assists the query, or
 - Helps validate, complete, or correct earlier reasoning.
- Ensure information is accurate, relevant, and objective.

Do NOT attempt to correct prior reasoning. Treat it as possibly flawed context.

Output Format:

****Final Information****

[Helpful factual information, or the non-knowledge-based response]

Inputs:

- Previous Reasoning Steps: {prev_reasoning}
- Current Search Query: {search_query}
- Wikipedia Content: {wikipedia_content}

Prompt for Summarizing Other Documents During Evaluation

Task Instruction:

You are tasked with reading and analyzing Wikipedia content based on the following inputs: **Previous Reasoning Steps**, **Current Search Query**, and **Wikipedia Content**. Your objective is to extract factual and relevant information from the **Wikipedia Content** that directly supports or informs the **Current Search Query**, and integrate it into the reasoning process in an objective and helpful manner.

Guidelines:

- Analyze Wikipedia Content:
 - Read carefully.
 - Identify factual info directly related to the query.
- Maintain Objectivity:
 - Do not validate or revise prior reasoning.
 - Use it as flawed context.

Output Format:

****Final Information****

[Helpful information]

Inputs:

- Previous Reasoning Steps: {prev_reasoning}
- Current Search Query: {search_query}
- Wikipedia Content: {wikipedia_content}

Prompt for Summarizing Other Documents During Training

Task Instruction:

Your first task is to determine whether the provided query is a knowledge-based query that can be answered using factual information from Wikipedia, or if it requires design, computation, or complex reasoning.

Step 1: Query Classification

- If knowledge-based (e.g., facts, definitions, history), proceed to Step 2.
- Otherwise, return:

This query requires design, computation, or complex reasoning, which exceeds the capabilities of a search engine. Please input another query or proceed with direct reasoning.

Step 2: Analyze Knowledge-Based Queries

- Read Wikipedia content;
- Extract relevant factual information;
- Stay neutral—do not alter previous reasoning;

Output Format:

****Final Information****

[Helpful information or the non-knowledge-based response]

Inputs:

- Previous Reasoning Steps: {prev_reasoning}
 - Current Search Query: {search_query}
 - Wikipedia Content: {wikipedia_content}
-

D.4 Prompts for Baseline Methods

Self-Ask Initial Prompt

Instruction:

The self-ask method uses few-shot examples to demonstrate the reasoning pattern:

Example 1:

Question: Who lived longer, Muhammad Ali or Alan Turing?

Are follow up questions needed here: Yes.

Follow up: How old was Muhammad Ali when he died?

Intermediate answer: Muhammad Ali was 74 years old when he died.

Follow up: How old was Alan Turing when he died?

Intermediate answer: Alan Turing was 41 years old when he died.

So the final answer is: Muhammad Ali

Example 2:

Question: When was the founder of craigslist born?

Are follow up questions needed here: Yes.

Follow up: Who was the founder of craigslist?

Intermediate answer: Craigslist was founded by Craig Newmark.

Follow up: When was Craig Newmark born?

Intermediate answer: Craig Newmark was born on December 6, 1952.

So the final answer is: December 6, 1952

Example 3:

Question: Who was the maternal grandfather of George Washington?

Are follow up questions needed here: Yes.

Follow up: Who was the mother of George Washington?

Intermediate answer: The mother of George Washington was Mary Ball Washington.

Follow up: Who was the father of Mary Ball Washington?

Intermediate answer: The father of Mary Ball Washington was Joseph Ball.

So the final answer is: Joseph Ball

Input:

Question: {question}

Options: {options}

Are follow up questions needed here:

Self-Ask Sub-question Answering Prompt

Instruction:

Please answer the following question based on the reference text. If the reference text does not contain sufficient information to answer the question, you may use your own knowledge to provide the answer. Always think step by step. Provide your final answer in the format `\boxed{YOUR_ANSWER}`.

Input:

- Question: {subquestion}
 - Reference text: {reference}
-

RAT Draft Generation Prompt

System Prompt:

You are an advanced AI assistant tasked with answering open-domain questions. You excel at providing comprehensive, well-structured answers with multiple paragraphs. Each paragraph you write contains multiple sentences that thoroughly explore the topic. You always follow formatting instructions precisely.

Instruction:

IMPORTANT: Structure your response as follows:

1. Write a comprehensive answer with **MULTIPLE PARAGRAPHS** (3-6 paragraphs typically).
2. Each paragraph **MUST** contain **AT LEAST 2** complete sentences. Single-sentence paragraphs are **NOT** acceptable.
3. Separate paragraphs with blank lines (press Enter twice).
4. At the very end, after all paragraphs, add your final answer in this format:

`\box{ANSWER}`

where ANSWER is **ONLY** the direct answer - typically just a name, number, date, or short phrase. Examples:

- For “Who was the first president?” → `\box{George Washington}`
- For “When was the company founded?” → `\box{1812}`
- For “What is the capital?” → `\box{Paris}`

DO NOT include explanations or full sentences in the box.

Input:

- Question: {question}

RAT Query Generation Prompt

Instruction:

Based on the question and the current answer content, generate a search query to verify or find additional information. Please summarize the content with the corresponding question. This summarization will be used as a query to search with Bing search engine. The query should be short but need to be specific to promise Bing can find related knowledge or pages. You can also use search syntax to make the query short and clear enough for the search engine to find relevant language data. Try to make the query as relevant as possible to the last few sentences in the content.

IMPORTANT: Just output the query directly. DO NOT add additional explanations or introduction in the answer unless you are asked to.

Input:

- Question: {question}
- Current Answer: {current_answer}

RAT Answer Revision Prompt

Instruction:

I want to revise the answer according to retrieved related text of the question. You need to check whether the answer is correct. If you find some errors in the answer, revise the answer to make it better. If you find some necessary details are ignored, add it to make the answer more plausible according to the related text.

IMPORTANT:

1. Keep the structure with multiple substantial paragraphs.
2. Use blank lines to separate paragraphs (press Enter twice).
3. If the original answer has `\box{ . . . }` at the end, you MUST keep it and update it if needed.
4. The `\box{ }` should contain ONLY the direct answer (name/number/date/short phrase), NOT a full sentence.

Just output the revised paragraphs directly, including the `\box{ }` if present.

Input:

- Retrieved Text: {retrieved_text}
- Question: {question}
- Answer: {current_answer}

Search-o1 Reasoning Prompt

System Prompt:

You are a reasoning assistant with the ability to perform web searches to help you answer the user's question accurately. You have special tools:

- To perform a search: write `<|begin_search_query|>` your query here `<|end_search_query|>`.
- Then, the system will search and analyze relevant web pages, then provide you with helpful information in the format `<|begin_search_result|>` ...search results... `<|end_search_result|>`.

You can repeat the search process multiple times if necessary. The maximum number of search attempts is limited to {max_rounds}.

Once you have all the information you need, continue your reasoning.

Example:

Question: "Alice David is the voice of Lara Croft in a video game developed by which company?"

Assistant thinking steps:

- I need to find out who voices Lara Croft in the video game.
- Then, I need to determine which company developed that video game.

Prompt for MMLU-Pro&MedQA (CoT)

Instruction:

You are solving a multiple-choice question. Think step by step and use careful reasoning. For each question, **analyze all options one by one**. For each option:

- Consider its meaning and implications.
- Evaluate whether it is correct or incorrect, and **explain why**.
- Eliminate incorrect options with clear, logical reasoning.

After analyzing all options, compare the remaining ones and choose the best answer.

At the end of your reasoning, give your final answer in the following format:

the correct answer is: A, B, C, D, etc. (only the letter corresponding to the correct option).

Input:

- Question: {question}
 - Options: {options}
-

Prompt for Math (CoT)

Instruction:

Please answer the following math question. You should think step by step to solve it.

Provide your final answer in the format `\boxed{YOUR_ANSWER}`.

Input:

- Question: {question}
-

Prompt for Open-Domain QA (CoT)

Instruction:

{question}

Please reason step by step, and put your final answer within `\boxed{}`.

E Case Study

E.1 Search independently based on different difficulty levels and questions

<p style="text-align: center; color: #00AEEF; font-weight: bold;">Simple questions are not searched</p> <p>Question: Let $z = 2 + \sqrt{2} - (3 + 3\sqrt{2})i$, and let $c = 2 - 3i$. Let w be the result when z is rotated around c by $\frac{\pi}{4}$ counter-clockwise. Ground Truth: $6 - 5i$</p>	<p style="text-align: center; color: #00AEEF; font-weight: bold;">Difficult questions are searched</p> <p>Question: In the state of Coinland, coins have values 6, 10, and 15 cents. Suppose x is the value in cents of the most expensive item in Coinland that cannot be purchased using these coins with exact change. What is the sum of the digits of x? Ground Truth: 11 . 0</p>
<p>Reasoning:</p> <p>To find the complex number w after rotating z around c by $\frac{\pi}{4}$ counter-clockwise, we will follow these steps:</p> <ol style="list-style-type: none"> Express z and c in complex form: $z = 2 + \sqrt{2} - (3 + 3\sqrt{2})i$ $c = 2 - 3i$ Translate z and c so that c is at the origin: $z' = z - c = (2 + \sqrt{2} - (3 + 3\sqrt{2})i) - (2 - 3i) = \sqrt{2} - 3\sqrt{2}i$ Rotate z' by $\frac{\pi}{4}$ counter-clockwise: The rotation of a complex number z' by θ counter-clockwise is given by multiplying z' by $e^{i\theta}$. For $\theta = \frac{\pi}{4}$: $e^{i\frac{\pi}{4}} = \cos\left(\frac{\pi}{4}\right) + i \sin\left(\frac{\pi}{4}\right) = \frac{\sqrt{2}}{2} + i \frac{\sqrt{2}}{2}$ <p>Therefore,</p> $z'' = z' \cdot e^{i\frac{\pi}{4}} = (\sqrt{2} - 3\sqrt{2}i) \left(\frac{\sqrt{2}}{2} + i \frac{\sqrt{2}}{2} \right)$ Perform the multiplication: $z'' = \sqrt{2} \cdot \frac{\sqrt{2}}{2} + \sqrt{2} \cdot i \frac{\sqrt{2}}{2} - 3\sqrt{2}i \cdot \frac{\sqrt{2}}{2} - 3\sqrt{2}i \cdot i \frac{\sqrt{2}}{2}$ <p>Simplify each term:</p> $z'' = \frac{2}{2} + i \frac{2}{2} - i \frac{6}{2} + \frac{6}{2}$ $z'' = 1 + i - 3i + 3$ <p>Combine like terms:</p> $z'' = 4 - 2i$ <div style="border: 1px solid red; padding: 2px; display: inline-block; margin-left: 100px;">No search</div> Translate z'' back by adding c: $w = z'' + c = (4 - 2i) + (2 - 3i) = 6 - 5i$ <p>Thus, the final answer is:</p> $6 - 5i$	<p>Reasoning:</p> <p style="color: green; font-weight: bold;">To solve this problem, we need to find the largest value that cannot be expressed as a combination of the coin values 6, 10, and 15 cents.</p> <p style="color: orange;"><se> Definition of the Frobenius coin problem for three coin values. </se> <info></p> <p style="color: green; font-weight: bold;">For three coin denominations, the Frobenius coin problem seeks the largest monetary amount that cannot be obtained using any combination of the given denominations. ...(more search information content) In the case of denominations 6, 10, and 15, since their GCD is 1, the Frobenius number is guaranteed to exist, but its exact value must be computed algorithmically or via specialized methods for three variables.</p> <div style="border: 1px solid green; padding: 2px; display: inline-block; margin-left: 100px;">Search</div> <p style="color: green; font-weight: bold;"></info></p> <p>Since there is no simple formula for three coin values, we will use a computational approach or known results to find the largest value that cannot be expressed as a combination of 6, 10, and 15.</p> <p>By checking manually or using a computational tool, we find that the largest value that cannot be expressed as a combination of 6, 10, and 15 is 29 cents.</p> <p>To confirm, we can check if 29 can be expressed as a combination of 6, 10, and 15:</p> <ul style="list-style-type: none"> • $29 - 15 = 14$ (not a combination of 6 and 10) • $29 - 10 = 19$ (not a combination of 6 and 10) • ... (more computational validation) <p>Since 29 cannot be expressed as a combination of 6, 10, and 15, we check the next value, 30, which can be expressed as:</p> $30 = 2 \times 15$ <p>Thus, the largest value that cannot be purchased using these coins is:</p> 29 <p>The sum of the digits of 29 is:</p> $2 + 9 = 11$

Figure 4: UR² can search independently based on different difficulty levels and questions. **Left:** For simpler questions, UR² considers that external knowledge support is not necessary and does not perform a search. **Right:** For questions of a certain difficulty, UR² will actively search for answers.

E.2 Planning

Planning	Planning
<p>Question: Hobbes describes felicity as:</p> <p>Options:</p> <ul style="list-style-type: none"> A. the absence of desire. B. a state of constant dissatisfaction. C. a state of constant change. D. the absence of progress. E. a state of satisfaction with one's current state. F. the fulfilment of all desires. G. a continual progress of desire from one object to another. H. the ability to control one's desires. I. a state of tranquility. J. a state of contemplation. <p>Ground Truth: G</p>	<p>Question: A 39-year-old woman is brought to the emergency department because of fevers, chills, and left lower quadrant pain. Her temperature is 39.1°C (102.3°F), pulse is 126/min, respirations are 28/min, and blood pressure is 80/50 mm Hg. There is blood oozing around the site of a peripheral intravenous line. Pelvic examination shows mucopurulent discharge from the cervical os and left adnexal tenderness. Laboratory studies show:</p> <ul style="list-style-type: none"> • Platelet count 14,200/mm³ • Fibrinogen 83 mg/mL (N = 200–430 mg/dL) • D-dimer 965 ng/mL (N < 500 ng/mL) <p>When phenol is applied to a sample of the patient's blood at 90°C, a phosphorylated N-acetylglucosamine dimer with 6 fatty acids attached to a polysaccharide side chain is identified. A blood culture is most likely to show which of the following?</p> <p>Ground Truth:D. Lactose-fermenting, gram-negative rods forming pink colonies on MacConkey agar</p>
<p>Reasoning: To determine which option best describes Hobbes's concept of felicity, let's analyze each option step by step. First, let's issue a query to understand Hobbes's definition of felicity.</p> <p><code><se></code> How does Hobbes define felicity? <code></se></code> <code><info></code></p> <div style="display: flex; justify-content: space-around;"> <div style="border: 1px solid black; padding: 2px;">Step 1</div> <div style="border: 1px solid black; padding: 2px;">Two step plan</div> </div> <p>According to Thomas Hobbes, "felicity" refers to a continual process of satisfying desires. ...<i>(more search information content)</i> This understanding is rooted in his broader view of human nature, where individuals are motivated by a continuous striving for things that they perceive as good for their wellbeing.</p> <p><code></info></code></p> <p>Based on the information provided, Hobbes describes felicity as a continual process of satisfying desires, which is best described by the following options:</p> <ul style="list-style-type: none"> A. the absence of desire. B. a state of constant dissatisfaction. C. a state of constant change. D. the absence of progress. E. a state of satisfaction with one's current state. F. the fulfilment of all desires. G. a continual progress of desire from one object to another. H. the ability to control one's desires. I. a state of tranquility. J. a state of contemplation. <div style="border: 1px solid black; padding: 2px; display: inline-block;">Step 2</div> <p>Thus, the correct answer is: G</p>	<p>Reasoning: Let's analyze each option step by step.</p> <div style="border: 1px solid black; padding: 2px; display: inline-block;">Task Decomposition</div> <ul style="list-style-type: none"> • The patient is experiencing fevers, chills, and left lower quadrant pain, which are consistent with a bacterial infection, likely pelvic inflammatory disease (PID) or a peritonitis. • ...<i>(more gradual thinking process)</i> • We need to determine which gram-negative bacteria are likely to be present in a blood culture. <p><code><se></code> Which gram-negative bacteria form colonies on blood agar? <code></se></code> <code><info></code></p> <div style="border: 1px solid black; padding: 2px; display: inline-block;">Information Seeking</div> <p>Several gram-negative bacteria are capable of forming colonies on blood agar, though their growth characteristics may vary: ...<i>(more search information content)</i> In summary, many gram-negative bacteria, especially members of the <i>Enterobacteriaceae</i> family, form colonies on blood agar. However, some gram-negative species (such as <i>Haemophilus influenzae</i>) require special conditions or enriched media for growth.</p> <p><code></info></code></p> <p>Based on the information provided, the patient is likely experiencing a gram-negative bacterial infection, possibly <i>Neisseria gonorrhoeae</i> or <i>Neisseria meningitidis</i>, given the characteristics of the colonies on blood agar.</p> <p>Given the options:</p> <div style="border: 1px solid black; padding: 2px; display: inline-block;">Evidence Integration</div> <ul style="list-style-type: none"> A. Coagulase-positive, gram-positive cocci forming mauve-colored colonies on methicillin-containing agar <p>...<i>(more analysis content)</i> Given the gram-negative nature of the bacteria and the typical growth on blood agar, the correct answer is: D.</p>

Figure 5: UR² can formulate plans and dynamically adjust them during the reasoning process. **Left:** UR² develops a plan and completes it in two steps. The first step is to search for the required knowledge, and the second step is to check each option individually. **Right:** UR² demonstrates clear planning behavior by decomposing the diagnostic task into sequential reasoning steps, identifying knowledge gaps, and querying external information to support its final decision.

E.3 Cross Validation

Cross Validation

Question:
What is the name of the 1999 film directed by Roman Polanski and starring his wife Emmanuelle Seigner that was a world-wide success but critically and commercially unsuccessful in North America?
Ground Truth: The Ninth Gate

Reasoning:
To answer this question, I need to identify the 1999 film directed by Roman Polanski that starred his wife Emmanuelle Seigner, and then determine if it was a worldwide success but critically and commercially unsuccessful in North America.
<se>
Who starred in films directed by Roman Polanski in 1999?
</se>
<info>
According to the Wikipedia content, Roman Polanski directed the film "The Ninth Gate" in 1999. The content does not provide a full cast list, but it confirms that "The Ninth Gate" is a 1999 mystery thriller film directed by Roman Polanski. **...(more search information content)** but it does confirm that "The Ninth Gate" is the relevant 1999 film directed by Roman Polanski.
</info>
Since "The Ninth Gate" was the 1999 film directed by Roman Polanski, I need to check if Emmanuelle Seigner starred in it and if it had the described success pattern.
<se>
Did Emmanuelle Seigner star in "The Ninth Gate"?
</se>
<info>
According to the Wikipedia content, Emmanuelle Seigner is known for her role in "The Ninth Gate" (1999). This confirms that Emmanuelle Seigner starred in "The Ninth Gate," the 1999 film directed by Roman Polanski.
</info>
Now that we know "The Ninth Gate" was the 1999 film directed by Roman Polanski and starring Emmanuelle Seigner, I need to verify if it was a worldwide success but critically and commercially unsuccessful in North America.
<se>
Was "The Ninth Gate" a worldwide success but critically and commercially unsuccessful in North America?
</se>
<info>
According to the Wikipedia content **...(more search information content)** However, despite its poor reception in North America, the film earned a worldwide gross of \$58.4 million against a \$38 million budget, indicating that it achieved commercial success internationally.
</info>
Based on the evidence, "The Ninth Gate" was indeed a worldwide success but critically and commercially unsuccessful in North America.
Thus, the correct answer is: The Ninth Gate

Layered Verification

Figure 6: This case demonstrates UR² cross validation ability by verifying the same candidate answer "The Ninth Gate" from multiple angles. It confirms the director and year, validates the actress's participation, and contrasts international vs. North American commercial and critical success. The multi-step confirmation strengthens the confidence in the final answer, showcasing robust cross validation behavior.

E.4 Reflection(Retry Retrieve/Refine Query)

Reflection (Retry Retrieve/Refine Query)

Question:

A 3-year-old boy is brought to the physician for the evaluation of recurrent skin lesions. The episodes of lesions started at the age of 2 months and multiple treatment options have been attempted without success. He has also had several episodes of respiratory tract infections, enlarged lymph nodes, and recurrent fevers since birth. The boy attends daycare. His older brother has asthma. The patient's immunizations are up-to-date. He is at the 5th percentile for length and 10th percentile for weight. He appears ill. His temperature is 38°C (100.4°F), pulse is 100/min, and blood pressure is 100/60 mm Hg. Examination shows several raised, erythematous lesions of different sizes over the face, neck, groins, and extremities; some secrete pus. Cervical and axillary lymph nodes are enlarged bilaterally. The remainder of the examination shows no abnormalities. Which of the following is the most likely diagnosis?

Options:

- A. Chronic granulomatous disease
- B. Ataxia-telangiectasia
- C. Atopic dermatitis
- D. Wiskott-Aldrich syndrome
- E. Chediak-Higashi syndrome

Ground Truth: A

Reasoning:

Let's analyze each option step by step.

- The patient has recurrent skin lesions, respiratory tract infections, enlarged lymph nodes, and recurrent fevers since birth.
- The patient is at the 5th percentile for length and 10th percentile for weight, appearing ill.
- The temperature is 38°C (100.4°F), pulse is 100/min, and blood pressure is 100/60 mm Hg.
- Examination shows raised, erythematous lesions of different sizes over the face, neck, groins, and extremities, some secrete pus. Cervical and axillary lymph nodes are enlarged bilaterally.

Let's issue a query to determine the likely diagnosis.

<se>

Which disease causes recurrent skin lesions and respiratory infections in a child?

</se>

<info>

Several rare immunodeficiency disorders can cause recurrent skin lesions and respiratory infections in children:

- **DOCK8 deficiency**: This is a combined immunodeficiency syndrome characterized by recurrent skin infections (bacterial, viral, and fungal), severe eczema, and recurrent respiratory tract infections such as sinusitis, otitis media, pneumonia, and bronchitis. Patients may also have recurrent fevers, lymphadenopathy (enlarged lymph nodes), and are prone to both skin and pulmonary infections.

...(more search information content)

These conditions are characterized by a combination of recurrent skin and respiratory infections, often accompanied by other signs of immune dysfunction such as lymphadenopathy, eczema, and growth failure. Among these, DOCK8 deficiency is particularly notable for the combination of severe skin infections, eczema, and recurrent respiratory tract infections in children.

</info>

Let's issue another query to narrow down the options.

<se>

Which immunodeficiency causes raised, erythematous lesions with pus?

</se>

<info>

Refine Query

Among primary immunodeficiencies, several can present with recurrent skin infections and lesions. According to the Wikipedia content, primary immune deficiency diseases often lead to increased susceptibility to recurrent skin infections, and some, such as Hyperimmunoglobulin E syndrome (also known as Job's syndrome), are specifically noted for causing skin lesions. ...(more search information content) erythematous, pus-filled skin lesions.

Therefore, Hyperimmunoglobulin E syndrome (Job's syndrome) is the immunodeficiency most characteristically associated with raised, erythematous skin lesions that may secrete pus.

</info>

Based on the information provided, the patient's symptoms are consistent with Hyperimmunoglobulin E syndrome (Job's syndrome), which is characterized by recurrent skin infections with pus-secreting lesions.

Thus, the correct answer is: A

Figure 7: After an initial analysis of the patient's symptoms, UR² conducted two rounds of retrieval: first a broad query ("Which disease causes..."), then refined to a targeted question ("Which immunodeficiency causes...") after recognizing the initial query was insufficient. This fully demonstrates the LLM's capability to refine queries.

Reflection (Retry Retrieve/Refine Query)

Question:

Marcia Baron argues that the artificiality of ticking bomb cases is problematic because

Options:

- A. the artificial assumptions are, in fact, empirical claims that are central to the debate.
- B. it leads us to believe that torture is always permissible.
- C. all artificial thought experiments are philosophically useless.

Ground Truth: A

Reasoning:

To determine which option correctly describes Marcia Baron's argument about the artificiality of ticking bomb cases, let's analyze the statement and each option step by step.

First, let's search for the specific reason Marcia Baron gives for why the artificiality of ticking bomb cases is problematic.

<se>

Why does Marcia Baron argue that ticking bomb cases are problematic?

</se>

<info>

The Wikipedia content provided does not include a direct summary of Marcia Baron's specific arguments about the problematic nature of ticking bomb cases. However...**(more search information content)** These hidden and often implausible assumptions obscure the true costs and complexities of permitting torture in real-life situations. While Marcia Baron's own detailed reasoning is not quoted in the provided content, it is common among philosophers critical of the ticking bomb scenario to highlight its artificiality and the way it distorts moral and practical considerations about torture.

</info>

Based on the information provided, it seems that Marcia Baron's argument likely focuses on the unrealistic assumptions inherent in ticking bomb cases. Let's compare this with the options given in the question.

Option A:

"the artificial assumptions are, in fact, empirical claims that are central to the debate."

This option seems to suggest that the artificial assumptions are actually empirical claims, which is not a common interpretation of Marcia Baron's argument. The argument is more about the unrealistic nature of the assumptions rather than their empirical status.

Option B:

"it leads us to believe that torture is always permissible."

This option seems to suggest that the artificiality leads to a false belief about the permissibility of torture, which is not a common interpretation of Marcia Baron's argument. The argument is more about the unrealistic assumptions rather than the permissibility of torture.

Option C:

"all artificial thought experiments are philosophically useless."

This option seems to suggest that all artificial thought experiments are useless, which is not a common interpretation of Marcia Baron's argument. The argument is more about the specific unrealistic assumptions in ticking bomb cases rather than the usefulness of all artificial thought experiments.

From the information provided, it seems that none of the options fit perfectly, but Option A might be the closest if we interpret "empirical claims" as a way to describe the assumptions, even though it's not the primary focus.

Let's refine with a more specific search:

<se>

What specific unrealistic assumptions does Marcia Baron criticize in ticking bomb cases?

</se>

<info>

The Wikipedia content does not provide a direct quotation or detailed summary of Marcia Baron's own words regarding the specific unrealistic assumptions she criticizes in ticking bomb cases. However, it does outline the general criticisms made by opponents of the ticking time bomb scenario—criticisms that are commonly attributed to philosophers like Baron. ...**(more search information content)** While Marcia Baron's specific phrasing is not included in the Wikipedia content, her critique is generally understood to align with these points about the unrealistic and misleading nature of the assumptions underlying ticking bomb cases.

</info>

Based on the detailed information provided, Option A seems to be the best fit, even though it might be a bit stretched:

Thus, the correct answer is: A

Refined Query

Figure 8: In this case, UR² initially retrieved general critiques related to ticking bomb scenarios but found no direct reference to Marcia Baron's argument. Recognizing this gap, it engaged in Reflection by refining the query to focus on "specific unrealistic assumptions", thereby exemplifying its capability to self-correct and deepen its information retrieval process.

Reflection (Retry Retrieve/Refine Query)

Question:

This question refers to the following information.

Source 1:

“The Sovereign Congress of Venezuela, to which authority the people of the Republic of the New Granada have voluntarily stood by. Considering:

1. That united in a single Republic, the provinces of Venezuela and the New Granada have all proportions and ways to elevate themselves to a higher grade of power and prosperity.

...(more source information content)

Source 2:

...(more source information content)

Question: Which of the following would most undermine the expectations expressed in Article 1 of Source 1?

Options:

- A. Intervention by the United States and European powers in the independence of Latin American countries
- B. Nationalism in Latin America pressing for stronger unity between the former Spanish colonies of South America
- C. Slave rebellions fragmenting newly independent Latin American states
- D. Sectionalism of elites in the various territories of newly independent Latin American states

Ground Truth: D

Reasoning:

To determine which option most undermines the expectations expressed in Article 1 of Source 1, let’s analyze the article and the options step by step.

Article 1 of Source 1 states:

“*The Republics of Venezuela and New Granada are from this day on united in a single one under the glorious title of Republic of Colombia.*”

This article expresses the expectation that the Republics of Venezuela and New Granada should be united to elevate themselves to a higher grade of power and prosperity.

Let’s search for a fact to help us determine which option most undermines this expectation.

<se>

What was the outcome of the union between Venezuela and New Granada?

</se>

<info>

The union between Venezuela and New Granada, along with Ecuador and Panama, resulted in the creation of the Republic of Gran Colombia in 1819. ...(more search information content) Thus, the intended lasting unity and prosperity were undermined by the eventual fragmentation of the union.

</info>

Based on the information provided, the correct option that most undermines the expectations expressed in Article 1 of Source 1 is the one that relates to the fragmentation of the union.

Option A:

“*Intervention by the United States and European powers in the independence of Latin American countries*”

This might have influenced the political landscape but does not specifically address the fragmentation of the union.

...(more analysis content)

This is more likely to be a factor in the fragmentation of the union, but let’s check the fact again.

Let’s refine:

<se>

What internal conflicts led to the dissolution of Gran Colombia?

</se>

<info>

The dissolution of Gran Colombia in 1830–1831 was primarily caused by internal political conflicts and regional tensions. The main internal conflicts included: ...(more search information content)

In summary, the dissolution of Gran Colombia was driven by internal political divisions over federalism versus centralism, strong regional identities and rivalries, and leadership struggles, all of which led to civil conflict and the eventual fragmentation of the union.

</info>

Based on the information provided, option D seems to be the best fit.

Thus, the correct answer is: D

Refined Query

Figure 9: UR² initially retrieved general information about Gran Colombia’s dissolution, but upon reflection recognized this was insufficient for precise answer selection. It then refined the query to investigate specific internal conflicts, ultimately identifying elite sectionalism as the key undermining factor, demonstrating effective use of Reflection through query refinement.

Reflection (Retry Retrieve/Refine Query)

Question:

Nathanson supports which of the following kinds of retributivism?

Options:

- A. proportional
- B. equality
- C. retaliatory
- D. punitive
- E. compensatory
- F. absolute
- G. restorative
- H. corrective
- I. preventive
- J. consequential

Ground Truth: A

Reasoning:

To determine which kind of retributivism Nathanson supports, let's analyze the concept of retributivism and the options provided.

First, let's search for specific information about Nathanson's stance on retributivism.

<se>

What kind of retributivism does Nathanson support?

</se>

<info>

The provided Wikipedia content does not mention Nathanson or specify what kind of retributivism he supports. However, it does define retributive justice as a theory of punishment that holds offenders should suffer in proportion to their offense, emphasizing proportionality and procedural standards, and distinguishing retribution from revenge. Without direct reference to Nathanson, no specific information about his stance or the particular kind of retributivism he supports can be extracted from the current Wikipedia content.

</info>

Based on the information provided, it seems we need more specific details about Nathanson's stance. Let's try another approach by looking at the general types of retributivism and see if any of them fit the description of Nathanson's support.

<se>

What are the main types of retributivism in criminal justice?

</se>

<info>

The Wikipedia content indicates that retributivism is a theory of punishment in criminal justice which holds that offenders deserve to suffer in proportion to their offenses. While the content does not list explicit "types" of retributivism by name, it does reference several key distinctions and purposes within retributive theories:

- **Proportionality**: Retributive justice emphasizes that the punishment should be proportional to the crime committed.
- **Desert-Based (Backward-Looking) Retributivism**: This form holds that offenders deserve punishment simply because they have committed a wrongdoing, focusing on blameworthiness and restoring balance (e.g., Lex Talionis or the Unfair Advantage Principle).
- **Expressive or Denunciatory Retributivism**: This type emphasizes the role of punishment in openly denouncing crime and expressing societal condemnation.
- **Restorative Elements**: Some retributive theories aim to restore balance or provide satisfaction to victims or society.

Additionally, the content distinguishes retributive justice from revenge by noting that retribution is not personal, has inherent limits, and follows procedural standards.

In summary, the main types or approaches within retributivism in criminal justice, as reflected in the Wikipedia content, include:

- Proportional (or "just deserts") retributivism
- Expressive/denunciatory retributivism
- Restorative or balance-restoring retributivism

</info>

Based on the information provided, let's compare the options with the types of retributivism mentioned:

Option A, proportional, seems to be the best fit based on the detailed description provided.

Thus, the correct answer is: A

Retry & Refine Query

Figure 10: This case demonstrates UR² capacity to reflect and refine its query strategy when initial retrieval fails. By shifting from a direct search for Nathanson's view to analyzing types of retributivism, the model successfully deduced the answer through indirect reasoning. This exemplifies its performance in reflection and retry-based retrieval refinement.