DEMYSTIFYING SCIENTIFIC PROBLEM-SOLVING IN LLMs by Probing Knowledge and Reasoning

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026027028

029

031

032

033

034

035

037

040

041

042

043

044

046

047

048

049

051

052

Paper under double-blind review

ABSTRACT

Scientific problem solving poses unique challenges for LLMs, requiring both deep domain knowledge and the ability to apply such knowledge through complex reasoning. While automated scientific reasoners hold great promise for assisting human scientists, there is currently no widely adopted holistic benchmark for evaluating scientific reasoning, and few approaches systematically disentangle the distinct roles of knowledge and reasoning in these tasks. To address these gaps, we introduce SCIREAS, a diverse suite of existing benchmarks for scientific reasoning tasks, and SCIREAS-PRO, a selective subset that requires more complex reasoning. Our holistic evaluation surfaces insights about scientific reasoning performance that remain hidden when relying on individual benchmarks alone. We then propose KRUX, a probing framework for studying the distinct roles of reasoning and knowledge in scientific tasks. Combining the two, we conduct an in-depth analysis that yields several key findings: (1) Retrieving task-relevant knowledge from model parameters is a critical bottleneck for LLMs in scientific reasoning; (2) Reasoning models consistently benefit from external knowledge added in-context on top of the reasoning enhancement; (3) Enhancing verbalized reasoning improves LLMs' ability to surface task-relevant knowledge.¹

1 Introduction

Recent frontier reasoning models, such as OpenAI's o-series (OpenAI et al., 2024) and DeepSeek-R1 (DeepSeek-AI et al., 2025), demonstrate significant advances by leveraging increased test-time compute to enable intermediate reasoning steps (Wei et al., 2023; Kojima et al., 2023). These approaches facilitate advanced mechanisms, including methodology exploration (Yao et al., 2023), self-verification (Ma et al., 2025a), and backtracking (Yang et al., 2025), resulting in improvements on tasks such as mathematics and coding with more test-time compute (Muennighoff et al., 2025).

These advances in reasoning capabilities create opportunities for applying LLMs to complex scientific tasks (Lu et al., 2024; Gottweis et al., 2025; Schmidgall et al., 2025). However, scientific work demands not only rigorous reasoning but also deep domain knowledge, from specialized concepts and foundational theories to hands-on methodological expertise and familiarity with obscure yet pivotal findings. Successful scientific reasoning systems must apply such knowledge in complex multi-step reasoning processes (Zhao et al., 2023; Wang et al., 2023a; Wadden et al., 2024a; Li et al., 2025).

While a variety of scientific benchmarks exist (e.g., GPQA (Rein et al., 2024) and MMLU-Pro (Wang et al., 2024b)), there is no holistic and unified benchmark that comprehensively targets scientific reasoning. Existing individual benchmarks typically focus narrowly on specific domains, task formats, or skill types. For example, although GPQA is challenging, it focuses exclusively on multiple-choice questions within a limited range of domains. Furthermore, there is a lack of analytical tools that can isolate the distinct roles that reasoning and scientific knowledge play when performing sophisticated scientific tasks.

We introduce datasets and methods to facilitate the study of scientific problem solving. First, we present **SCIREAS**, a unified suite of ten public benchmarks that span physics, chemistry, biology, medicine, materials, mathematics, computer science, and engineering, with multiple-choice,

¹The codebase and artifacts are released at link-redacted-for-review.

056

058

060

061

062

063

064

065

066

067

068

069

071

073

074

075

076 077

078

079

081 082

084

085

087

090

092

094

095

096

098 099

100

101

102

103

104

105

106

107

Figure 1: KRUX pipeline. Starting upper left, we prompt an LLM (one of base, DeepSeek-R1, Base-Math, Base-STEM, and Base-BOTH) with a question from SCIREAS as the knowledge source, collect the output and reasoning traces, and feed the reasoning traces to DeepSeek-R1 as the extractor to generate knowledge ingredients (KIs). We then evaluate the tested model with KI-augmented questions, which allows us to study three key research questions (RQ1 §4.2.2, RQ2 §4.2.3, RQ3 §4.2.4) regarding LLMs' knowledge and reasoning capabilities in scientific problem-solving.

fill-in-the-blank, structured, and protocol/procedural questions. To improve evaluation efficiency and sharpen the focus on reasoning difficulty, we manually inspect each subtask and retain only those that are subject-relevant and reasoning-intensive, while preserving broad domain coverage. Furthermore, to facilitate standardized evaluation, we provide an efficient and unified implementation of streamlined assessment across individual benchmarks, avoiding the need to set up different environments or dataset-specific boilerplate for each dataset (§3.1).

Next, we introduce **SCIREAS-PRO**, a compact subset of SCIREAS tailored for evaluating more challenging reasoning. Specifically, SCIREAS-PRO is constructed by selecting examples from SCIREAS where only reasoning models with high inference-time compute budget (or the highest allowed number of thinking tokens) succeed. We find that despite containing only 8% as many examples as SCIREAS, SCIREAS-PRO better differentiates weak and strong reasoners (§3.1).

Having constructed the reasoning-intensive scientific benchmarks, our next goal is to leverage them to study how verbalized chain-of-thought (CoT) reasoning affects knowledge recall and usage (§4). To study this, we design **KRUX** (Knowledge & Reasoning Utilization eXams), a probing framework which supplies models with atomic "knowledge ingredients" (KIs) extracted from other models' reasoning traces (Figure 1). This technique allows for more controlled analyses of reasoning and knowledge, which we use to perform three in-depth investigations that lead to the following findings:

- (1) Vanilla instruct models can *outperform* their reasoning counterparts by $\geq 10\%$ once KIs are provided in-context, suggesting that internalizing and retrieving the right knowledge is a key bottleneck for scientific reasoning tasks.
- (2) When both model families receive the same KIs from a strong reasoner (e.g., DeepSeek-R1), the reasoning-fine-tuned models consistently outperform the base models, showing that **reasoning** models are capable of utilizing external in-context knowledge for additional improvements.
- (3) Feeding KIs from a reasoning-fine-tuned model to its base model can boost performance even when the KIs are already known by the base model, indicating that **reasoning-fine-tuning aids knowledge recall by surfacing more relevant knowledge**.

Our contributions can be summarized as:

- We introduce SCIREAS, a unified and holistic benchmark suite spanning a broad range of scientific domains and problem types, allowing us to surface insights that otherwise remain hidden if relying on individual datasets only. We also release a reasoning-focused subset SCIREAS-PRO that allows efficient benchmarking of sophisticated reasoning with more room for improvement.
- We present KRUX, a novel analytic framework which we use to conduct a comprehensive empirical study that disentangles the impacts of knowledge and reasoning.
- We provide an in-depth analysis with three key findings: (i) knowledge retrieval is a bottleneck; (ii) in-context knowledge consistently benefits reasoning models; and (iii) long CoT improves knowledge surfacing. We support these findings with controlled post-training experiments, and show our training recipe is competitive compared with concurrent SFT post-training efforts.

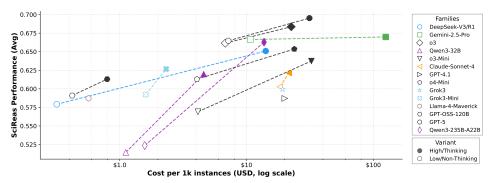


Figure 2: Frontier reasoning models' performance evaluated on SCIREAS. The X-axis shows the cost per 1k instances in USD. Different reasoning settings on the same model can result in distinct costs and performance, but the margins vary depending on the models.

2 RELATED WORK

Scientific Benchmarks Existing scientific benchmarks span a wide array of domains and tasks, but each tends to focus on specific disciplines or subskills, often lacking explicit emphasis on multistep reasoning or standardized implementation. For example, most tasks in SciRIFF (Wadden et al., 2024a) focus on context-grounded information QA, rather than demanding reasoning. Benchmarks like GPQA (Rein et al., 2024) and LabBench (Laurent et al., 2024) pose reasoning challenges, yet they cover only a limited range of scientific domains and rely on multiple-choice QA formats. Implementation-wise, benchmarks lack standardized prompts, up-to-date evaluation metrics, or consistent scoring and reporting, making reproducibility and fair comparison difficult (Gu et al., 2025; Gao et al., 2024). To address this fragmentation, our study systematically incorporates 10 prominent scientific benchmarks, GPQA, MMLU-Pro (Wang et al., 2024b), SuperGPQA (Team et al., 2025b), LabBench, OlympiadBench (He et al., 2024), SciBench (Wang et al., 2023b), SciRIFF, UG-Physics (Xu et al., 2025), SciEval (Sun et al., 2024), and SciKnowEval (Feng et al., 2024), enabling a unified, comprehensive, and reproducible evaluation suite of scientific reasoning capabilities.

Knowledge & Reasoning An important line of work on disentangling reasoning and knowledge designs specialized tasks (e.g., linguistically challenging questions (Bean et al., 2024; Khouja et al., 2025) or synthetic multi-hop questions (Li & Goyal, 2025)) to isolate reasoning from knowledge, but such benchmarks are often artificial and domain-constrained. Notably, Li & Goyal (2025) analyzes the synergy between knowledge and reasoning as knowledge evolves, offering a perspective complementary to our controlled CoT SFT experiments. Another line of work trains external classifiers to label questions as reasoning- or knowledge-intensive based on parametric models (Thapa et al., 2025). However, this approach requires well-calibrated training data and does not distinguish the tested model's internal knowledge from reasoning. Concurrent work leverages reasoning traces to evaluate factual correctness (Wu et al., 2025), but focuses on surface-level factuality rather than genuine knowledge recall. Unlike prior work that trains external classifiers to label question types or checks surface factuality in traces, KRUX holds knowledge constant and varies the target model, isolating knowledge recall from reasoning ability without relying on heuristic difficulty tags. Additional related work is provided in Appendix B.

3 BENCHMARKING KNOWLEDGE-INTENSIVE SCIENTIFIC REASONING

Given limited coverage in terms of domain, formats, or accessibility for individual benchmarks, SCIREAS solves this by merging ten datasets under one standardized harness, offering broad domain coverage and consistent evaluation.

3.1 EVALUATION SUITE CONSTRUCTION

SCIREAS SCIREAS is an evaluation suite focused on reasoning-intensive scientific tasks curated from 10 representative existing benchmarks. Through task-level filtering, SCIREAS reduces instance count by nearly 50% while preserving coverage, and, inspired by OLMES (Gu et al., 2025), provides

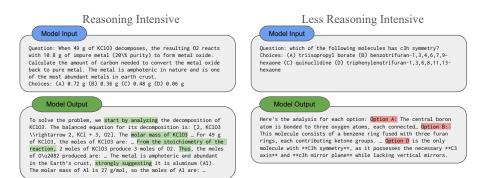


Figure 3: An example pair with varying reasoning intensity, where the example on the left is sampled from SCIREAS-PRO and the right is a filtered out example (§3.1). On the left, the progressive reasoning chain is highlighted. The example on the right emphasizes knowledge recall on each option with a simple elimination strategy.

a unified implementation optimized with vLLM (Kwon et al., 2023) and batch job APIs² for scalable, easy-to-use, and efficient evaluation.

Our curation prioritizes subtasks from each benchmark that demand not only specific domain knowledge but also complex, multi-step reasoning processes for resolution. For each subtask from each benchmark, we used a subtask-level exclusion protocol: for each candidate subtask, we sampled 20 instances and excluded the subtask if any sample failed to both (i) require domain knowledge beyond the prompt and (ii) require multi-step reasoning. This deliberately conservative, exclusion-oriented design looks for reasons to remove subtasks, biasing against inclusion and reducing the risk of false positives.³ We provide a complete list of selected subtasks in Appendix C.1

To keep evaluation cost-efficient, we uniformly sample 200 instances from each subtask sourced from high-cost benchmarks — MMLU-Pro, SciKnowEval, SciEval, and UGPhysics, which maintains similar evaluation outcomes (more in Appendix C.2) while reducing the cost by nearly 50% (from 29,604 to 15,567 total instances). Benchmarks affected by our filtering are marked with an asterisk (*); their scores are not directly comparable to those from prior work.

SCIREAS-PRO Although SCIREAS provides a uniform measurement for model performance on scientific reasoning tasks that nominally require scientific reasoning, the difficulty of individual instances is uneven: some can be answered with little deductive effort once the pertinent fact is recalled, as shown in an example in Figure 3.

To isolate the reasoning skill, we therefore curate a "hard" subset — those questions whose solutions still demand multi-step inference even when all relevant knowledge is available — so that any performance gains cannot be explained by knowledge recall alone. Building on our observation in §3.2, we hypothesize that the performance difference under different test-time inference budgets can serve as an effective indicator of reasoning intensity. Specifically, instances where reasoning models *fail* with low reasoning budget but *succeed* with high budget likely require complex reasoning, even when the necessary domain knowledge is accessible to the model in both settings.

In practice, we evaluate o3-mini and o4-mini on SCIREAS with both *high* and *low* "reasoning-effort" settings, an OpenAI API flag that limits the number of thinking tokens before output. For o3-mini and o4-mini, the high-effort setting costs about $5.8 \times$ more per instance than the low-effort setting (Table 6, Appendix C.1).⁴ For each model, we keep questions answered *incorrectly* under low effort but *correctly* under high effort and take the union of these sets to create SCIREAS-PRO, resulting in 1,260 unique instances. We further validate this approach by using LLM judge and human evaluation to check the reasoning-intensiveness of resulting examples from this filtering pipeline in Appendix C.3, and observe that incorrect answers are attributed to insufficient reasoning rather than lack of knowledge 90% of the time by humans on a sampled set and 91% of the time by LLM judge.

²We provide batch job inference options for popular LLM providers, e.g., OpenAI, Anthropic, TogetherAI, and Gemini. Using batch APIs allows for up to 50% cost reduction.

³While this manual inspection can be subjective, it is based on the authors' graduate-school-level expertise.

⁴Because these models are proprietary, factors beyond the flag may influence performance. We therefore treat the flag as a practical, not absolute, proxy and validate it with independent studies (Appendix C.3).

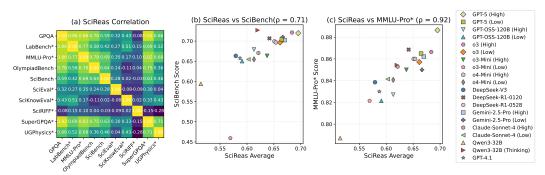


Figure 4: SCIREAS correlations breakdown. (a) Task-to-task Pearson correlations. SCIREAS incorporates tasks complementary to popular benchmarks. (b) and (c) show performance on SCIREAS vs. SciBench and MMLU-Pro*. Models may be tuned for certain tasks, outperforming higher-ranked models on individual benchmarks.

3.2 BENCHMARKING FRONTIER MODELS

Having constructed SCIREAS and SCIREAS-PRO with focus on scientific reasoning tasks, we now examine how state-of-the-art models perform under varying computational budgets. We evaluate frontier models using different "reasoning-effort" settings (see configuration details in Appendix D). These settings typically correspond to significant differences in output length, with high-effort modes producing substantially more reasoning tokens as they work through complex problems.⁵

Aggregated Results Figure 2 highlights aggregated performance evaluated on SCIREAS, with score breakdowns on selected models shown in Table 6. Notably, the **aggregated ranking provides additional insights that differ from popular individual benchmarks**. Comparing o3-High and Gemini-2.5-Pro-Preview-High as an example, o3-High wins on GPQA and MMLU-Pro* while Gemini-2.5-Pro-Preview-High wins on SuperGPQA*, all with a thin margin (within 1 absolute point, even evaluated on MMLU-Pro before uniform sampling as shown in Figure 7). Similarly, GPT-5-High shows on-par performance with Gemini-2.5-Pro-Preview-High on problem-solving benchmarks like OlympiadBench and SciRIFF. Evaluated across SCIREAS, however, we notice that GPT-5-High outperforms its competitors on a broader range of benchmarks. Meanwhile, o3-High achieves higher overall performance over Gemini-2.5-Pro-Preview-High, with superior performance on LabBench* and weaker on OlympiadBench by a large margin (beyond 10 absolute points).

Benchmark Correlations In general, as the Pearson correlations shown in Figure 4 (a), while some benchmarks are closely correlated (e.g., GPQA and SuperGPQA*), benchmarks containing free-form QA and fill-in-the-blank questions like SciRIFF* and SciEval* are not highly correlated with GPQA-like multiple-choice tasks, demonstrating the need for a holistic evaluation suite. Isolating specific benchmarks, we observe that models from different providers may be tuned explicitly for specific tasks or skills. As shown in Figure 4 (b) and (c), Qwen3-32B-Thinking strikes noticeably above the trend on SciBench, reaching comparable performance to commercial frontier models. Similarly, DeepSeek-V3 and DeepSeek-R1-0120 demonstrate stronger performance on MMLU-Pro*, indicating capabilities that surpass their overall rankings.

Performance Gap by Reasoning Difference Although the gap varies depending on different model families, **the same model can exhibit a significant performance gap under different reasoning settings**. For instance, in Figure 2, o3-mini-Low and -High show a performance gap of 6.8. Similar traits can be observed among o4-mini, Claude-Sonnet-4, and o3, while Gemini-2.5-Pro-Preview shows the least performance gain, even with significantly more (>10×) thinking budget. This observation motivates the construction of SCIREAS-PRO, leveraging the performance gap between low and high reasoning efforts as an effective proxy for identifying instances that demand complex reasoning rather than mere knowledge recall. **For practitioners, task-specific evaluation is still recommended** for the optimal balance between inference cost and performance.

Amplified Performance Gap Figure 5 shows that **SCIREAS-PRO amplifies performance gaps between low- and high-reasoning settings**, where the gap between GPT-5-High and GPT-5-Low

⁵In this work, we refer to DeepSeek-R1-0528 and DeepSeek-V3-0324 simply as DeepSeek-R1 and DeepSeek-V3, respectively, unless otherwise specified.

widens from 3.01 to 12.22, and the corresponding gap for Gemini-2.5-Pro-Preview widens from 0.35 to 2.30. Meanwhile, non-reasoning models, e.g., GPT-4.1, DeepSeek-V3, show more significant gaps compared to concurrent reasoning models, o3 and DeepSeek-R1, respectively.

Reasoning Efforts Improve Math Reasoning More Is a higher inference budget more helpful to math or numeric reasoning than non-math reasoning? To answer this question, we categorize instances from SCIREAS into *Has-Math* and *No-Math* buckets (Appendix E.3.1) and report the gains in micro average accuracy. In Appendix

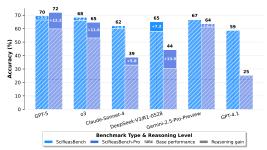


Figure 5: Model performance on SCIREAS and SCIREAS-PRO with varying reasoning capabilities. SCIREAS-PRO amplifies gaps between low-reasoning and high-reasoning settings.

E.3.2 Figure 10, the results show that higher reasoning budgets yield more improvements among Has-Math instances compared to No-Math instances. This finding echoes with concurrent work where Sprague et al. (2024) points out that CoT helps more with math and symbolic reasoning.

4 DISENTANGLING KNOWLEDGE AND REASONING IN SCIENTIFIC TASKS

While SCIREAS and SCIREAS-PRO provide benchmarks to evaluate scientific reasoning capabilities, another fundamental question remains: how does CoT reasoning adaptation affect a model's ability to recall and utilize knowledge? To address this question, we first conduct a series of controlled SFT experiments on high-quality reasoning traces with and without in-domain scientific knowledge, and then we propose KRUX, a novel investigative framework to study three key research questions regarding the role of knowledge in scientific reasoning using the fine-tuned checkpoints.

4.1 CONTROLLED COT SFT

To control for data composition and isolate the impact of reasoning and knowledge injection during post-training, we fine-tune Qwen2.5-7B-

Table 1: Performance of reasoning models trained from Qwen2.5-Instruct and Llama-3.1-Instruct on SYNTHETIC-1 and concurrent reasoning models.

Model	Method	SCIREAS	-Pro						
Our Checkpoints									
Qwen	_	37.07	13.97						
Qwen-STEM	SFT	40.47	16.11						
Qwen-Math	SFT	41.99	18.17						
Qwen-BOTH	SFT	42.84	21.11						
Llama	_	31.25	11.67						
Llama-STEM	SFT	35.28	14.29						
Llama-Math	SFT	35.49	16.98						
Llama-BOTH	SFT	38.55	16.51						
Concurrent I	Reasoning P	ost-training							
SYNTHETIC-1-SFT	SFT	37.64	19.44						
OpenR1	SFT	43.08	26.43						
Llama-Nemotron	SFT&RL	43.53	23.75						
General-Reasoner	RL	34.99	13.73						

Instruct (Yang et al., 2024) and Llama-3.1-8B-Instruct (Grattafiori et al., 2024) on reasoning traces drawn from mathematics and STEM domains, as well as on their combination. This allows us to attribute behavior changes to the data mixture rather than confounding factors.

For training, we leverage the SYNTHETIC-1 (Mattern et al., 2025) dataset, an existing large-scale dataset released by Prime Intellect, which consists of outputs of DeepSeek-R1-0120, including the reasoning traces, on a diverse set of tasks. More specifically, we leverage the mathematics and STEM subsets from SYNTHETIC-1 (denoted as SYNTHETIC-1-Math/STEM, respectively). The former provides reasoning traces on abstract math reasoning questions, serving as a source for long CoT adaptation without introducing in-domain knowledge. In contrast, the latter is sourced from StackExchange (Lambert et al., 2023), providing a more in-domain data source for a broader range of scientific subjects. The math subset contains around 462K instances, while the STEM subset contains around 512K instances. Details of the training and evaluation setup are in Appendix E.

By training Qwen2.5-7B-Instruct on SYNTHETIC-1 (-Math, -STEM, and the combined subsets), we derived Qwen-Math, Qwen-STEM, and Qwen-BOTH along with their counterparts trained from Llama-3.1-8B-Instruct. In the following, we will refer to the Base models as Qwen or Llama for

⁶Notably, SYNTHETIC-1-Math is sourced from competition-level math problems, highlighting high-quality abstract math reasoning filtered by verified answers. In contrast, StackExchange and SYNTHETIC-1-STEM provide more realistic problem-solving data from wider subjects, offering more coverage in science domains.

brevity. Compared with concurrent work on long CoT post-training (Bercovich et al., 2025a; Face, 2025; Mattern et al., 2025; Ma et al., 2025b), our checkpoints deliver comparable performance under controlled settings (Table 1), serving as reliable investigating checkpoints. A lightweight analysis of domain-specific improvements and data composition is presented in Appendix E.3-E.4.

4.2 Knowledge & Reasoning Utilization Exam (KRUX)

We introduce KRUX (Figure 1), a novel investigative framework to study the role of knowledge and long CoT reasoning in scientific problem solving. To separate what a model *knows* from how it *reasons*, we hold knowledge availability fixed by injecting compact, answer-agnostic knowledge ingredients (KIs) in-context. In the framework, we extract KIs from the reasoning traces of various models and provide these KIs in-context to LLMs when evaluating them. Consequently, gains over a no-KI baseline indicate a knowledge bottleneck, while persistent errors point to reasoning limits.

We first introduce our pipeline to extract KIs from reasoning traces (§4.2.1), and then discuss how we analyze and apply extracted KIs to test knowledge recall (§4.2.2, §4.2.4) and usage (§4.2.3). For experiments, we prioritize challenging benchmarks (e.g., GPQA, MMLU-Pro*, and LabBench*), which have been widely used by previous work in the field on tasks that require scientific expertise.

4.2.1 Knowledge Ingredient (KI) Extraction

First, to analyze the role of knowledge in models' performance on scientific problem-solving, we aim to study a setting in which the model is given the requisite knowledge in-context. Specifically, we take the reasoning traces from a reasoning model as the knowledge source and use a strong reasoning-focused LLM (e.g., DeepSeek-R1) to extract the essential atomic knowledge units that comprise it, which we refer to as *knowledge ingredients* (KIs) (Figure 1). We provide the extraction prompt and example KIs in Appendix F.1. We then augment the original question by prepending the extracted KIs in-context and ask the models to solve the same problem.

We perform additional checks on DeepSeek-R1 and Qwen3-30B-A3B-Thinking-2507 as the extractor to ensure that KIs (a) are task-agnostic (i.e., provide knowledge and facts without referring to specific details in the question or options, e.g, "... as referred in option B ..."), (b) do *not* leak any part of the final answer, and (c) strictly adhere to traces as the knowledge source without additional information. In manual review, *all* extracted KIs from 100 sampled reasoning traces met these criteria and were consistent with their source reasoning traces. For the following analysis, we use KIs generated by DeepSeek-R1 with more details on the differences between extrators in Appendix F.2.

To prevent the extractor from hallucinating or introducing extraneous facts (i.e., KIs unsupported by the source trace or unnecessary for solving the problem), we feed the generated KIs back to the source model and measure performance. If performance changes materially, this indicates potential leakage of steps or answers. Empirically, we observe no significant change (Table 2, Base vs. w/ Base KIs), suggesting the KIs are answer-agnostic and faithful to the trace. Further, although it is possible that the knowledge pieces may be irrelevant to the solution, as shown in recent studies of CoT faithfulness (Turpin et al., 2023; Wang et al., 2024c;a), recent high-performing models like DeepSeek-R1 have demonstrated strong reasoning adherence on benchmark tasks (DeepSeek-AI et al., 2025). Our experiments show that the knowledge pieces help models on reasoning tasks. See Figures 13-15 in Appendix F.1 for KI examples generated by different models for the same question.

Centered on our main objective of studying knowledge recall and use in reasoning models, we examine the following key research questions: **RQ1:** To what extent can base models benefit from high-quality external knowledge? **RQ2:** Do reasoning-enhanced models benefit from external knowledge? **RQ3:** Does reasoning fine-tuning improve models' ability to surface helpful knowledge?

4.2.2 RQ1: TO WHAT EXTENT CAN BASE MODELS BENEFIT FROM HIGH-QUALITY EXTERNAL KNOWLEDGE?

Problem Statement. We investigate the potential improvement from external knowledge by providing KIs to the base models in the prompt when performing scientific reasoning (Figure 1). Here, we focus on two sources for the KIs, which are extracted from their own CoT traces (w/ Base KIs) or from DeepSeek-R1's CoT traces (w/ R1 KIs). To overcome context sensitivity, we report averages and standard deviations across 5 runs with corresponding KIs permuted randomly. We

then investigate whether there are significant gaps between base models augmented with additional KIs in the context, and their corresponding reasoning-fine-tuned models. To this end,

comparisons are made with reasoning-fine-tuned models trained on our controlled data mixtures and the ones from concurrent work (i.e., General-Reasoner-7B (Liu et al., 2025) and Llama-Nemotron-Nano-8B (Bercovich et al., 2025b)) that involve SFT and reinforcement learning based on the same base models.

Answer to RQ1: As an upper bound, a base model with high-quality in-context knowledge can substantially outperform its reasoning-enhanced counterpart.

As shown in Table 2, base models provided with KIs from DeepSeek-R1 are able to *outperform* base models alone or Base w/ Base KIs setup by $\geq 20\%$, and *outperform* reasoning variants without KIs by $\geq 10\%$ across different benchmarks and model families, showing the external knowledge provides greater gain than reasoning finetuning. The fact that a base model without strong reasoning capabilities can outperform reasoning models in this setting indicates a potential deficiency of the models in knowledge recall that hinders their performance in scientific reasoning.

Table 2: Performance on GPQA and LabBench* with base models alone, base models with KIs extracted from DeepSeek-R1 or itself (w/ {R1, Base} KIs), and reasoning-fine-tuned models. Best and second best average scores are labeled in bold and underlined. Reasoning models fall behind base models augmented with in-context knowledge.

Setup	GPQA	LabBench*
Qwen	35.27	32.38
w/ Qwen KIs	34.24 ± 0.93	30.93 ± 1.43
w/R1 KIs	47.19 ± 1.53	41.40 ± 2.46
Qwen-STEM	<u>41.63</u>	31.75
Qwen-Math	39.47	30.18
Qwen-BOTH	40.81	33.83
General-Reasoner	35.94	<u>35.58</u>
Llama	28.13	33.55
w/ Llama KIs	29.06 ± 1.44	34.40 ± 2.58
w/R1 KIs	43.57 ± 0.88	42.27 ± 1.60
Llama-STEM	38.95	36.04
Llama-Math	36.16	34.78
Llama-BOTH	<u>39.43</u>	<u>36.61</u>
Llama-Nemotron	37.95	27.78

4.2.3 RQ2: Do reasoning-enhanced models benefit from external knowledge?

Table 3: Accuracy of Qwen and Llama variants on benchmarks with external knowledge ingredients (KIs). We report averages and standard deviations over 5 random permutations of the KIs. Reasoning variants w/ R1 KIs outperform base model w/ R1 KIs across different benchmarks and models.

	GP	QA	MML	U-Pro*	LabBench*		
Models	w/ self KIs	w/ R1 KIs	w/ self KIs	w/ R1 KIs	w/ self KIs	w/R1 KIs	
Qwen	34.24 ± 0.93	47.19 ± 1.53	59.03 ± 0.34	68.86 ± 0.56	30.93 ± 1.43	41.40 ± 2.46	
Qwen-STEM	41.63 ± 2.10	52.50 ± 2.14	64.71 ± 1.05	69.69 ± 0.73	31.75 ± 2.81	43.79 ± 1.71	
Qwen-Math	39.47 ± 1.66	53.53 ± 1.24	66.93 ± 0.72	74.00 ± 0.59	30.18 ± 1.65	41.17 ± 2.32	
Qwen-BOTH	40.81 ± 2.04	54.46 ± 1.27	65.71 ± 0.74	71.64 ± 1.16	33.83 ± 2.59	43.90 ± 2.71	
Llama	29.06 ± 1.44	43.57 ± 0.88	47.73 ± 0.89	60.53 ± 1.67	34.40 ± 2.58	42.27 ± 1.60	
Llama-STEM	38.95 ± 1.31	53.17 ± 1.15	59.14 ± 0.85	68.19 ± 1.15	36.04 ± 3.98	46.87 ± 1.49	
Llama-Math	36.16 ± 2.33	53.75 ± 1.15	59.65 ± 0.98	69.01 ± 0.55	34.78 ± 4.26	45.55 ± 0.68	
Llama-BOTH	39.43 ± 2.00	54.73 ± 1.75	63.81 ± 0.90	72.74 ± 0.26	36.61 ± 2.73	48.65 ± 0.49	

Problem Statement. Observing considerable improvements from adding external KIs from DeepSeek-R1 to base models in RQ1, we hypothesize similar improvements would scale on reasoning-enhanced models, offering additional gains on top of enhanced reasoning. To this end, we evaluate base and CoT SFTed variants on KIs extracted from DeepSeek-R1, providing the same necessary knowledge from DeepSeek-R1's reasoning traces (w/ R1 KIs). As a baseline without the added knowledge, we provide the tested models with KIs extracted from their own CoT traces (w/ self KIs) for comparison.

Answer to RQ2: Yes. Reasoning models also substantially benefit from addition of contextual knowledge. As shown in Table 3, within both Qwen and Llama groups, reasoning-enhanced models w/ R1 KIs in the context show significant improvements over the base setting without the KIs, while preserving the gap compared with the base model w/ R1 KIs. Confirming the effectiveness of providing external knowledge as an in-context prompt, this result sheds light on potential future improvement by applying high-quality external memory modules as an external knowledge source for better problem-solving capabilities, echoing the finding in COMPACTDB (Lyu et al., 2025), a concurrent effort constructing a high-quality datastore for reasoning-intensive tasks.

Note, however, that in these experiments, we do not distinguish between two possible non-exclusive explanations for the improvement from adding R1 KIs. (a) It may be that the R1 KIs provide new key knowledge absent from the model's parameters, or (b) the model may already possess these facts but struggle to retrieve them (put another way, once a strong reasoning model supplies the *key* facts, the reasoning search space might narrow and the problem becomes easier, whether or not the model originally "knew" the augmented facts). We further analyze this confounder in RQ3.

4.2.4 RQ3: Does reasoning fine-tuning improve models' ability to surface helpful knowledge?

Problem Statement. While we observe that external knowledge benefits reasoning models, in this RQ, we ask how reasoning-fine-tuning affects knowledge recall. To this end, we focus on evaluating the KIs from -Math models to determine whether they offer more improvement than those of base models, as -Math models are fine-tuned on math-only data without additional scientific knowledge.

Notably, in Table 2, while -STEM and -BOTH variants, trained with SYNTHETIC-1-STEM, outperform -Math variants due to science in-domain training data, -Math variants also largely outperform the base model even without being trained on science data. Recalling our discussion in RQ2 (§4.2.3), the -Math model's gains have the same two non-exclusive explanations, (a) the -Math model performs better on science questions that require math

Table 4: Accuracy (%) of synthetic knowledge recall on KIs generated from Qwen/Llama-Math on GPQA and MMLU-Pro*. Base models and math reasoning-fine-tuned models show similar performance on knowledge recall questions.

KI Dataset	Qwen Qwen	-Math -Math	Llama	-Math -Math
KI-GPQA	72.30	73.02	70.94	68.94
KI-MMLU-Pro*	82.49	81.50	74.46	74.12

because math knowledge was loaded into the model through the math-specific fine-tuning, and/or (b) the -Math model is better at surfacing its relevant parametric knowledge via CoT expression.

To disentangle these two factors, we extract KIs from the CoTs of the -Math models and examine whether these KIs represent new knowledge added by fine-tuning, or whether they are also facts known to the base model. We probe this by querying the model with synthetic questions that test knowledge of each KI (see Appendix F.3 for examples). Then, to verify explanation (b), we provide the KIs in-context from either the -Math or base model, to the corresponding base model; i.e., holding mathematical reasoning capacity constant while varying only the external knowledge.

Answer to RQ3: Yes. In response to explanation (a), we find that on average, the base models and their corresponding -Math variants have similar recall of the KIs (Table 4), meaning that explanation (a) is unlikely to be the major contributor for the improvements.

To verify explanation (b), Table 5 shows that KIs from -Math deliver significant boosts over those from the base models across different benchmarks and model families. This result suggests that CoT verbalization improves the model's ability to surface the most relevant knowledge for the given reasoning problems. Notably, the KIs are unlikely to have been newly acquired during fine-tuning (Table 4); instead, the findings indicate that reasoning-fine-tuned models exhibit improved recall of knowledge already parameterized in the base model.

Table 5: Performance on GPQA and MMLU-Pro* with KIs extracted from base and -Math reasoning models. KIs extracted from -Math models enable more improvement over those from base models.

Base	Setup	GPQA	MMLU-Pro*
Qwen	w/ Qwen KIs w/ Qwen-Math KIs	34.24 ± 0.93 36.93 ± 1.75	59.03 ± 0.34 63.66 ± 0.45
Llama	w/ Llama KIs w/ Llama-Math KIs		47.73 ± 0.89 53.91 ± 0.94

5 CONCLUSION

In this work, we studied how reasoning and domain knowledge each contribute to scientific reasoning in LLMs. To this end, we introduced SCIREAS and SCIREAS-PRO, unified, reproducible suites for evaluating scientific reasoning across domains and formats, together with KRUX, a knowledge-controlled evaluation framework. We showed: (i) retrieving task-relevant knowledge from parameters is a key bottleneck; (ii) reasoning-fine-tuned models get complementary gains from external KIs; and (iii) verbalized CoT improves knowledge surfacing. Our results show that reasoning-focused fine-tuning improves both reasoning and knowledge use, suggesting promising future directions in better understanding and enhancing these interconnected components.

REFERENCES

486

487

488

489

490 491

492

493

494

495

496

497

498

499 500

501

502

504

505

506

507

509

510

511

512

513

514

515

516

517

519

520

521

522

523

524

525

526 527

528

529

530

531 532

533

534

536

- Iván Arcuschin, Jett Janiak, Robert Krzyzanowski, Senthooran Rajamanoharan, Neel Nanda, and Arthur Conmy. Chain-of-thought reasoning in the wild is not always faithful, 2025. URL https://arxiv.org/abs/2503.08679.
- Andrew Michael Bean, Simeon Hellsten, Harry Mayne, Jabez Magomere, Ethan A Chi, Ryan Andrew Chi, Scott A. Hale, and Hannah Rose Kirk. LINGOLY: A benchmark of olympiad-level linguistic reasoning puzzles in low resource and extinct languages. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=cLga8GStdk.
- Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, pp. 3615–3620. Association for Computational Linguistics, 2019.
- Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp, Ehud Karpas, Ran Zilberstein, Jiaqi Zeng, Soumye Singhal, Alexander Bukharin, Yian Zhang, Tugrul Konuk, Gerald Shen, Ameya Sunil Mahabaleshwarkar, Bilal Kartal, Yoshi Suhara, Olivier Delalleau, Zijia Chen, Zhilin Wang, David Mosallanezhad, Adi Renduchintala, Haifeng Qian, Dima Rekesh, Fei Jia, Somshubra Majumdar, Vahid Noroozi, Wasi Uddin Ahmad, Sean Narenthiran, Aleksander Ficek, Mehrzad Samadi, Jocelyn Huang, Siddhartha Jain, Igor Gitman, Ivan Moshkov, Wei Du, Shubham Toshniwal, George Armstrong, Branislav Kisacanin, Matvei Novikov, Daria Gitman, Evelina Bakhturina, Jane Polak Scowcroft, John Kamalu, Dan Su, Kezhi Kong, Markus Kliegl, Rabeeh Karimi, Ying Lin, Sanjeev Satheesh, Jupinder Parmar, Pritam Gundecha, Brandon Norick, Joseph Jennings, Shrimai Prabhumoye, Syeda Nahida Akter, Mostofa Patwary, Abhinav Khattar, Deepak Narayanan, Roger Waleffe, Jimmy Zhang, Bor-Yiing Su, Guyue Huang, Terry Kong, Parth Chadha, Sahil Jain, Christine Harvey, Elad Segal, Jining Huang, Sergey Kashirsky, Robert McQueen, Izzy Putterman, George Lam, Arun Venkatesan, Sherry Wu, Vinh Nguyen, Manoj Kilaru, Andrew Wang, Anna Warno, Abhilash Somasamudramath, Sandip Bhaskar, Maka Dong, Nave Assaf, Shahar Mor, Omer Ullman Argov, Scot Junkin, Oleksandr Romanenko, Pedro Larroy, Monika Katariya, Marco Rovinelli, Viji Balas, Nicholas Edelman, Anahita Bhiwandiwalla, Muthu Subramaniam, Smita Ithape, Karthik Ramamoorthy, Yuting Wu, Suguna Varshini Velury, Omri Almog, Joyjit Daw, Denys Fridman, Erick Galinkin, Michael Evans, Shaona Ghosh, Katherine Luna, Leon Derczynski, Nikki Pope, Eileen Long, Seth Schneider, Guillermo Siman, Tomasz Grzegorzek, Pablo Ribalta, Monika Katariya, Chris Alexiuk, Joey Conway, Trisha Saar, Ann Guan, Krzysztof Pawelec, Shyamala Prayaga, Oleksii Kuchaiev, Boris Ginsburg, Oluwatobi Olabiyi, Kari Briski, Jonathan Cohen, Bryan Catanzaro, Jonah Alben, Yonatan Geifman, and Eric Chung. Llama-nemotron: Efficient reasoning models, 2025a. URL https://arxiv.org/abs/2505.00949.
- Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, et al. Llama-nemotron: Efficient reasoning models. *arXiv preprint arXiv:2505.00949*, 2025b.
- Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. Crawling the internal knowledge-base of language models, 2023. URL https://arxiv.org/abs/2301.12810.
- DeepSeek-AI. Deepseek-r1: Usage recommendations, 2025. URL https://huggingface.co/deepseek-ai/DeepSeek-R1#usage-recommendations.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,

Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

- Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL https://github.com/huggingface/open-r1.
- Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. Sciknoweval: Evaluating multi-level scientific knowledge of large language models, 2024. URL https://arxiv.org/abs/2406.09098.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL https://zenodo.org/records/12608602.
- Daniela Gottesman and Mor Geva. Estimating knowledge in large language models without generating a single token, 2024. URL https://arxiv.org/abs/2406.12673.
- Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist, 2025. URL https://arxiv.org/abs/2502.18864.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Yuling Gu, Oyvind Tafjord, Bailey Kuehl, Dany Haddad, Jesse Dodge, and Hannaneh Hajishirzi. Olmes: A standard for language model evaluations, 2025. URL https://arxiv.org/abs/2406.08446.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. OlympiadBench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3828–3850, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.211. URL https://aclanthology.org/2024.acl-long.211/.

- Mingyu Jin, Weidi Luo, Sitao Cheng, Xinyi Wang, Wenyue Hua, Ruixiang Tang, William Yang
 Wang, and Yongfeng Zhang. Disentangling memory and reasoning ability in large language
 models, 2025. URL https://arxiv.org/abs/2411.13504.
 - Jude Khouja, Karolina Korgul, Simi Hellsten, Lingyi Yang, Vlad Neacsu, Harry Mayne, Ryan Kearns, Andrew Bean, and Adam Mahdi. Lingoly-too: Disentangling reasoning from knowledge with templatised orthographic obfuscation, 2025. URL https://arxiv.org/abs/2503.02972.
 - Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023. URL https://arxiv.org/abs/2205.11916.
 - Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
 - Nathan Lambert, Lewis Tunstall, Nazneen Rajani, and Tristan Thrush. Huggingface h4 stack exchange preference dataset, 2023. URL https://huggingface.co/datasets/HuggingFaceH4/stack-exchange-preferences.
 - Jon M Laurent, Joseph D Janizek, Michael Ruzo, Michaela M Hinks, Michael J Hammerling, Siddharth Narayanan, Manvitha Ponnapati, Andrew D White, and Samuel G Rodriques. Lab-bench: Measuring capabilities of language models for biology research. *arXiv preprint arXiv:2407.10362*, 2024.
 - Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
 - Aochong Oliver Li and Tanya Goyal. Memorization vs. reasoning: Updating Ilms with new knowledge, 2025. URL https://arxiv.org/abs/2504.12523.
 - Sihang Li, Jin Huang, Jiaxi Zhuang, Yaorui Shi, Xiaochen Cai, Mingjun Xu, Xiang Wang, Linfeng Zhang, Guolin Ke, and Hengxing Cai. Scilitllm: How to adapt llms for scientific literature understanding, 2025. URL https://arxiv.org/abs/2408.15545.
 - Qianchu Liu, Sheng Zhang, Guanghui Qin, Timothy Ossowski, Yu Gu, Ying Jin, Sid Kiblawi, Sam Preston, Mu Wei, Paul Vozila, Tristan Naumann, and Hoifung Poon. X-reasoner: Towards generalizable reasoning across modalities and domains, 2025. URL https://arxiv.org/abs/2505.03981.
 - Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery, 2024. URL https://arxiv.org/abs/2408.06292.
 - Xinxi Lyu, Michael Duan, Rulin Shao, Pang Wei Koh, and Sewon Min. Frustratingly simple retrieval improves challenging, reasoning-intensive benchmarks, 2025. URL https://arxiv.org/abs/2507.01297.
 - Ruotian Ma, Peisong Wang, Cheng Liu, Xingyan Liu, Jiaqi Chen, Bang Zhang, Xin Zhou, Nan Du, and Jia Li. S²r: Teaching llms to self-verify and self-correct via reinforcement learning, 2025a. URL https://arxiv.org/abs/2502.12853.
 - Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhu Chen. General-reasoner: Advancing Ilm reasoning across all domains. 2025b. URL https://api.semanticscholar.org/CorpusID:278768680.
 - Justus Mattern, Felix Gabriel, and Johannes Hagemann. Synthetic-1 release: Two million collaboratively generated reasoning traces from deepseek-r1, February 2025. URL https://www.primeintellect.ai/blog/synthetic-1-release.

649

650

651

652

653

654

655

656

657

658

659

661

662

663

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

684

685

686

687

688

689

690

691

692

693

696

697

699

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL https://arxiv.org/abs/2501.19393.

OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai ol system card, 2024. URL https://arxiv.org/abs/2412.16720.

Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window extension of large language models, 2023. URL https://arxiv.org/abs/2309.00071.

Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. Kilt: a benchmark for knowledge intensive language tasks, 2021. URL https://arxiv.org/abs/2009.02252.

Arvind Prabhakar et al. Omniscience: A domain-specialized llm for scientific reasoning. *arXiv* preprint arXiv:2503.17604, 2025.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a

- benchmark. In *First Conference on Language Modeling (COLM)*, 2024. URL https://openreview.net/forum?id=Ti67584b98.
- Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants, 2025. URL https://arxiv.org/abs/2501.04227.
- Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning. *ArXiv*, abs/2409.12183, 2024. URL https://api.semanticscholar.org/CorpusID:272708032.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. Scieval: A multi-level large language model evaluation benchmark for scientific research, 2024. URL https://arxiv.org/abs/2308.13149.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025a.
- P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixin Deng, Shawn Gavin, Shian Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, David Ma, Yuansheng Ni, Haoran Que, Qiyao Wang, Zhoufutu Wen, Siwei Wu, Tyshawn Hsing, Ming Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang, Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin, Yinghao Ma, Tianyang Pang, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu, Shanghaoran Quan, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jinyang Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Oian Liu, Tianyu Liu, Shiwen Ni, Junran Peng, Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and Ge Zhang. Supergpqa: Scaling Ilm evaluation across 285 graduate disciplines, 2025b. URL https://arxiv.org/abs/2502.14739.
- Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
- Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL https://qwenlm.github.io/blog/qwq-32b/.
- Rahul Thapa, Qingyang Wu, Kevin Wu, Harrison Zhang, Angela Zhang, Eric Wu, Haotian Ye, Suhana Bedi, Nevin Aresh, Joseph Boen, Shriya Reddy, Ben Athiwaratkun, Shuaiwen Leon Song, and James Zou. Disentangling reasoning and knowledge in medical large language models. *ArXiv*, abs/2505.11462, 2025. URL https://api.semanticscholar.org/CorpusID:278714970.
- Andrew Turpin, Jason Wei, Denny Zhou, Quoc V Le, and Ed H Chi. Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2305.15020*, 2023.
- David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle Lo, Tom Hope, Luca Soldaini, Shannon Zejiang Shen, Doug Downey, Hannaneh Hajishirzi, and Arman Cohan. Sciriff: A resource to enhance language model instruction-following over scientific literature, 2024a. URL https://arxiv.org/abs/2406.07835.
- David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle Lo, Tom Hope, Luca Soldaini, Shannon Zejiang Shen, et al. Sciriff: A resource to enhance language model instruction-following over scientific literature. *arXiv preprint arXiv:2406.07835*, 2024b.

- Changyue Wang, Weihang Su, Qingyao Ai, Yujia Zhou, and Yiqun Liu. Decoupling reasoning and knowledge injection for in-context knowledge editing, 2025. URL https://arxiv.org/abs/2506.00536.
- Pengfei Wang et al. Scienceqa: A large-scale open dataset for question answering in science education. *arXiv preprint arXiv*:2210.08127, 2023a.
 - Weijie Wang, Xiang Chen, et al. Evaluating the faithfulness of chain-of-thought reasoning in large language models. *arXiv preprint arXiv:2401.02392*, 2024a.
 - Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023b.
 - Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. arXiv preprint arXiv:2406.01574, 2024b. URL https://arxiv.org/abs/2406.01574.
 - Yunfan Wang, Dian Yu, Qian Zhou, et al. Can large language models follow chain-of-thought prompts faithfully? In *International Conference on Learning Representations (ICLR)*, 2024c.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
 - Juncheng Wu, Sheng Liu, Haoqin Tu, Hang Yu, Xiaoke Huang, James Zou, Cihang Xie, and Yuyin Zhou. Knowledge or reasoning? a close look at how llms think across domains, 2025. URL https://arxiv.org/abs/2506.02126.
 - Xin Xu, Qiyun Xu, Tong Xiao, Tianhao Chen, Yuchen Yan, Jiaxin Zhang, Shizhe Diao, Can Yang, and Yang Wang. Ugphysics: A comprehensive benchmark for undergraduate physics reasoning with large language models, 2025. URL https://arxiv.org/abs/2502.00334.
 - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
 - Xiao-Wen Yang, Xuan-Yi Zhu, Wen-Da Wei, Ding-Chu Zhang, Jie-Jing Shao, Zhi Zhou, Lan-Zhe Guo, and Yu-Feng Li. Step back to leap forward: Self-backtracking for boosting reasoning of language models, 2025.
 - Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.
 - Ge Zhang et al. Sciglm: Pre-training generalist language models for science with scientific papers. arXiv preprint arXiv:2402.00730, 2024.
 - Wayne Xin Zhao et al. A survey of llms for scientific research. arXiv preprint arXiv:2307.07927, 2023.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL https://arxiv.org/abs/2306.05685.

A LIMITATIONS

Our KRUX framework and KI extraction methods depend on strong models like DeepSeek-R1 for generating reasoning traces. While we used an open-weight model, which provides more transparency and interpretability, the KI extraction pipeline may introduce unobservable biases (though risk is minimal due to our focus on scientific domains), unwanted leakage of information about the answer, or inconsistencies in the faithfulness of the KIs to the task. To mitigate this, we conducted manual analysis of the KIs, confirming their relevance and no direct answer leakage, but extracted KIs could occasionally be irrelevant or incomplete, especially if deployed at scale. Furthermore, some of our analyses are confounded by factors such as context sensitivity (addressed via permutations) and the impact of constraining the search space when providing KIs, which we interpret as an upper bound but may overestimate pure recall benefits. We have taken measures to mitigate these and discussed the caveats in our discussion of results with more details.

Our experiments focus on moderate-sized LLMs with ¡10B parameters, specifically open-weight models (Qwen2.5, Llama3.1). While we deliberately selected two model families and models large enough to exhibit non-trivial reasoning performance, this limits the generalizability of our findings to larger models. Experimenting with larger models represents a straightforward extension but requires significantly greater computational resources, beyond the scope of our current study and our available compute resources.

The benchmarks we examine emphasize STEM fields, which may underrepresent interdisciplinary or emerging scientific research areas. We acknowledge potential data contamination issues that may impact our analysis; however, the nature of our study is analytical, and we perform controlled experiments. In our benchmarks, we also mitigate these concerns by focusing on recent 2024–2025 datasets. Despite these constraints, our methodology provides a systematic framework for evaluating domain-specific reasoning that can be extended to address these limitations in future work.

B EXTENDED RELATED WORK

Evaluating Knowledge of LLMs Early efforts tended to evaluate the LM knowledge frontier with a static unified benchmark (Petroni et al., 2021). However, given the growing training corpus for pushing LLM performance, quantifying the knowledge frontier of LLMs becomes increasingly challenging, making it difficult to design a unified benchmark. Instead of general knowledge evaluation, recent work approaches the knowledge frontier of LLMs by anchoring on specific entities, proposing methods to quantify knowledge and factuality around given entities (Gottesman & Geva, 2024; Cohen et al., 2023). With recent development of reasoning LLMs, more work exploits long CoT traces as evidence of explicit knowledge utilization, verifying knowledge recall in CoT traces for factuality (Wu et al., 2025). Nevertheless, directly evaluating CoT traces can result in false positive signals on the knowledge boundary, given that the knowledge involved could be factual but not helpful for problem solving (Arcuschin et al., 2025). In our framework, we construct controlled settings and protocols to evaluate whether the knowledge is genuinely helpful for problem-solving, implicitly guaranteeing the factuality and relevance.

Reasoning LLMs Recent work has shown that LLMs can be trained to utilize intermediate tokens for reasoning, achieving better performance on reasoning tasks as the decoding budget increases. OpenAI's o-series (OpenAI et al., 2024) represents the landmark of this paradigm among commercial frontier models, followed by DeepSeek-R1 (DeepSeek-AI et al., 2025) and several recent efforts to reproduce this success without releasing the training data, such as QwQ (Team, 2025) and Kimi (Team et al., 2025a). Some recent initiatives aim to achieve the same goal using fully open data sources, led by Llama-Nemotron from NVIDIA (Bercovich et al., 2025b) and SYNTHETIC-1 from Prime Intellect (Mattern et al., 2025), releasing post-training data to foster development within the community. Our work builds on these commitments, sharing the vision of improving model reasoning by leveraging intermediate tokens, while emphasizing our focus on scientific domains rather than on mathematics or general logical reasoning.

LLMs for Science Recent advancements in scientific LLMs have transitioned from early domain-specific pretraining (e.g., Beltagy et al. 2019; Lee et al. 2020), to more comprehensive models with multiple stages of training, e.g., SciGLM (Zhang et al., 2024), SciLitLLM (Li et al., 2025), and

OmniScience (Prabhakar et al., 2025). On the other hand, reasoning models have shown strong performance on scientific tasks such as GPQA and MMLU-Pro (DeepSeek-AI et al., 2025; OpenAI et al., 2024), and some recent efforts instrument LLMs to separate recall from deduction during inference (Wang et al., 2025; Jin et al., 2025). However, we still lack a clear understanding of the factors underlying performance on scientific tasks, such as knowledge acquisition or improved reasoning capabilities. We aim to address this gap by studying these factors and then providing a recipe for training more capable models in science.

C SCIREAS DETAILS

Benchmark		о3			o3-mii	ni	(04-min	ni	Ger	nini-2.	5-Pro	Clau	de-So	nnet-4	(GPT-5	5
	Low	High	Δ	Low	High	Δ	Low	High	Δ	Low	High	Δ	Low	High	Δ	Low	High	Δ
GPQA	75.4	79.9	+4.5	63.4	73.9	+10.5	69.4	74.6	+5.2	80.1	79.5	-0.6	63.8	69.0	+5.2	79.2	82.4	+3.1
SuperGPQA*	54.9	59.5	+4.6	40.5	54.0	+13.5	48.6	57.1	+8.5	60.1	60.4	+0.3	45.2	49.8	+4.6	58.6	62.4	+3.8
MMLU-Pro*	85.7	86.6	+0.9	82.1	85.0	+2.9	84.1	86.0	+1.9	85.0	86.2	+1.2	84.1	85.3	+1.2	86.5	88.6	+2.1
LabBench*	70.5	74.2	+3.7	56.9	59.2	+2.3	59.7	63.7	+4.0	61.9	64.4	+2.5	53.4	57.2	+3.8	66.6	74.4	+7.8
OlympBench	53.5	58.0	+4.5	39.5	51.1	+11.6	40.4	49.6	+9.2	67.5	69.6	+2.1	55.4	59.8	+4.4	60.0	64.9	+4.8
SciBench	69.7	72.1	+2.4	46.0	66.3	+20.3	65.5	69.7	+4.2	71.0	70.2	-0.8	65.5	67.1	+1.6	70.4	72.0	+1.6
SciEval*	84.8	82.7	-2.1	83.8	83.4	-0.4	87.1	87.5	+0.4	86.4	85.1	-1.3	85.8	85.8	0.0	87.4	86.1	-1.3
SciKnowEval*	52.1	51.9	-0.2	49.0	51.9	+2.9	49.9	51.1	+1.2	46.8	47.6	+0.8	43.6	43.3	-0.3	45.5	46.7	+1.2
SciRIFF*	51.8	53.6	+1.8	51.3	51.8	+0.5	50.6	52.2	+1.6	51.6	51.4	-0.2	53.5	50.9	-2.6	46.9	50.1	+3.3
UGPhysics*	63.1	65.2	+2.1	56.7	60.7	+4.0	57.7	62.2	+4.5	56.0	55.4	-0.6	52.4	53.2	+0.8	63.6	67.6	+4.0
Average	66.2	68.4	+2.2	56.9	63.7	+6.8	61.3	65.4	+4.1	66.6	67.0	+0.4	60.3	62.1	+1.8	66.5	69.5	+3.1
0.01\$ / Instance	0.68	2.25	$\times 3.3$	0.41	3.24	$\times 7.9$	0.41	2.38	$\times 5.8$	1.07	12.51	$\times 11.7$	1.83	7.50	$\times 4.1$	0.72	3.10	$\times 4.3$

Table 6: Performance (%) across SCIREAS grouped by models at low and high reasoning efforts. The same model with different reasoning effort can have distinctive performance with a clear margin.

C.1 EVALUATION SUITE CURATION

See Table 11-12 for domain distribution. We list the selection of each benchmark as follows.

GPQA (Rein et al., 2024): No change. Report in micro average. License: CC-BY-4.0.

MMLU-Pro (Wang et al., 2024b): MMLU-Pro features subjects beyond STEM and scientific subjects. We first filter by subjects, retaining instances from physics, chemistry, computer science, math, biology, and health, and then randomly sample each task to 200 instances max. Report in macro average across 7 subjects. License: MIT.

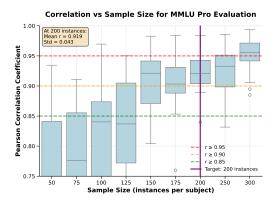
LabBench (**Laurent et al., 2024**): We drop tasks that require visual inputs or external table/paper extraction, therefore dropping DbQA, FigQA, LitQA2, SuppQA, and TableQA, retaining CloningScenarios, PropotolQA, and SeqQA. Report in macro average across 3 tasks. License: CC-BY-SA-4.0.

SciBench (Wang et al., 2023b): No change. Report in micro average. License: MIT.

OlympiadBench (**He et al., 2024**): Dropping tasks that require visual inputs or not in English. Report the macro average across math and physics. License: apache-2.0.

SciRIFF (Wadden et al., 2024b): We drop tasks that primarily focus on information/relation/table extraction and retain EvidenceInference, Qasper, and SciFact. Report in macro average of 5 metrics (detailed in Table 11-12) across 3 tasks. License: ODC-BY.

SciKnowEval (**Feng et al., 2024**): The authors introduce scientific tasks in 5 progressive levels from knowledge memorization to application. After manual inspection, we only preserve tasks from the highest level of knowledge application (L5), and cap instances from each task to be 200. Report the macro average across 8 tasks. License: MIT.



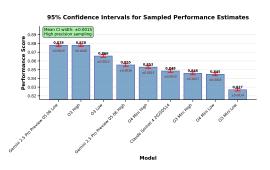


Figure 6: Correlation between sampled and full dataset performance as a function of sample size. 200 instances per subject (purple) yields $r = 0.919 \pm 0.043$. Error bars: SD over 30 samples.

Figure 7: 95% confidence intervals for 200-instance sampling across nine frontier models. Mean CI half-width ≈ 0.0015 ; numbers above/below bars show mean and half-width.

SciEval (Sun et al., 2024): Similar to SciKnowEval, the authors introduce 4 progressive levels of static tasks, including basic knowledge, knowledge application, scientific calculation, and research creativity. After inspection, we retain knowledge application and scientific calculation subsets, capping each task to a maximum of 200. Report the macro average across 6 tasks. License: N/A.

UGPhysics (**Xu et al., 2025**): The authors annotate each instance into 5 different physics reasoning skills: knowledge recall, laws application, math derivation, practical application, and others, which fails to be categorized into the categories above. We filter out instances that specifically require knowledge recall only and cap instances from each subject to be 200 max. Report the macro average across 13 subjects. License: CC-BY-NC-SA-4.0.

SuperGPQA (**Team et al., 2025b**): We curate questions from two broad domains — science and engineering — while omitting niche areas that lie outside mainstream STEM (e.g., weapon science, textile engineering). The science portion spans mathematics, biology, physics, systems science, and chemistry. The engineering portion covers a comprehensive set of disciplines: electronic science and technology; nuclear science and technology; mechanical engineering; information and communication engineering; civil engineering; instrument science and technology; computer science and technology; control science and engineering; chemical engineering and technology; mechanics; electrical engineering; materials science and engineering; hydraulic engineering; power engineering and engineering thermophysics; and optical engineering. Report in macro average across the domain of science and engineering. License: ODC-BY.

C.2 UNIFORM SAMPLING VALIDATION: MMLU-PRO CASE STUDY

Evaluating state-of-the-art frontier models could be expensive. To mitigate evaluation cost, we evaluate frontier models on MMLU-Pro* before and after uniform sampling. By sample size correlation in Figure 6 and 95% confidence intervals for sampled subset in Figure 7, we show that the sampling is cost-efficient and statistically effective while reducing evaluating instances from 6,696 to 1,400.

For costly frontier reasoning models such as Gemini-2.5-Pro-Preview, at rates in time of writing, the sampling reduces SCIREAS evaluation costs from \$3,600 to \$1,500 and can be further decreased to \$730 by using batch job inference.

C.3 Scireas-Pro Reasoning Intensiveness Validation

To test this hypothesis, we pursue two complementary checks: (1) different reasoning models should have high agreement identifying reasoning intensive instances, and (2) filtered instances should agree with human judgment in terms of reasoning intensiveness.

C.3.1 Cross-Model Agreement on Reasoning Intensity

To validate our hypothesis that performance gaps between different reasoning effort settings indicate reasoning intensity, we first examine whether different models agree on which instances are reasoning-intensive. As shown in Figure 1, for each reasoning model, we categorize each test question from SCIREAS by their correctness under low/high reasoning efforts into four categories, (high_c, low_c), (high_c, low_i), (high_i, low_c), and (high_i, low_i), where high/low stands for high/low reasoning effort setting and *_c/*_i stands for the problem instance has been answered correctly/incorrectly by the model. Treating (high_c, low_i) as targeting instances that require high reasoning effort, we measure how (high_c, low_i) sets derived from different reasoning models agree with others.

As shown in Table 7, treating (high_c, low_i) from o3-mini as ground truth, the same set derived from o4-mini, o3, and claude-sonnet-4 largely coincide with o3-mini across different benchmarks from SCIREAS (all above 70%), showing high agreement on instances that require high reasoning efforts across models from different model families.

Ground Truth Target	o3-mini vs. o4-mini	o3-mini vs. o3	o3-mini vs. claude-sonnet-4
SuperGPQA*	78.0	77.8	76.1
GPQA	80.4	81.0	79.0
MMLU-Pro*	92.2	91.6	92.0
LabBench*	71.9	74.6	75.8
SciBench	75.9	74.1	75.4
OlympiadBench	81.1	81.5	81.2
SciEval*	94.3	93.1	93.5
UGPhysics*	83.2	82.9	83.8

C.3.2 HUMAN AND LLM-AS-JUDGE ASSESSMENT

Table 7: Accuracy of overlapping instances on (high_c, low_i) from o3-mini vs. other models, treating o3-mini as ground true label. Different reasoning models agree on high reasoning instances.

The overlap of instances that require high reasoning effort shows reasoning models tend to agree on problem difficulty, but to verify the reliability of reasoning effort as a surrogate, the filter should also align with human judgment.

To this end, we collect the union of (highle, low_i) from o3-mini and o4-mini for the case study and apply an LLM-as-judge assessment (Zheng et al., 2023) to expedite the process while manually annotating a subset for a reliability test. The LLM judge is based on GPT-4.1 for a balanced tradeoff between assessment reliability and cost. Notably, naively prompting the LLM judge to determine the reasoning difficulty could be suboptimal due to a lack of reference. Therefore, we designed two reference-based evaluation protocols: (a) pair-wise comparison on reasoning difficulty between instance questions sampled from filtered subset and original SCIREAS, and (b) identifying failing reason for filtered instances given low and high reasoning outputs (i.e., whether the model fails in a low reasoning setting due to lack of reasoning effort).

(a) Pairwise Comparison For each instance in SCIREAS-PRO, the judge is also presented with an instance drawn from the set of other, non-overlapping instances from SCIREAS. The judge is not given any information as to which instance is drawn from which source and is tasked to identify which instance is more reasoning-intensive.

(b) Failure Analysis For each instance in SCIREAS-PRO, the judge is presented with both the correct high reasoning output (if both o3-mini-high and o4-mini-high are correct, o4-mini-high will be selected) as well as the incorrect low reasoning output from the corresponding model (e.g. correct: o3-mini-high; incorrect: o3-mini-low). The judge is tasked with determining whether the failure of the low reasoning effort model can be attributed primarily due to insufficient reasoning ability or lack of domain knowledge.

SYSTEM MESSAGE

1026

1027 1028

1029

1030

1031

1032

1033 1034

1035

1036

1039

1040

1041

1042

1043

1045

1046

1047

1048 1049 1050

1051 1052 1053

1054

1056

1058

1062 1063

1064

1067

1068

1069

1070

1071

1074

1077

1078

1079

You are an expert judge comparing reasoning intensity between two questions. Analyze both questions thoroughly and determine which one demands more complex reasoning.

Reply in this exact format:

```
\#\#\#EXPLANATION\colon <\!\! detailed analysis of both questions and the comparison> \#\#\#RESULTS\colon A / B / UNCLEAR
```

USER MESSAGE

You will be shown two questions (A and B) from the same academic domain.

A question is *reasoning intensive* if it requires:

- · Complex multi-step logical reasoning
- · Advanced mathematical computation or derivation
- · Integration of multiple concepts or principles
- Abstract thinking or sophisticated problem-solving strategies
- Deep domain knowledge application

```
*QUESTION A*
```

```
Context: {{context_a}}
Question: {{question_a}}
*QUESTION B*
Context: {{context_b}}
Question: {{question_b}}
```

Analyze both questions carefully and explain your reasoning. Then reply using the exact format speci-

fied above.

Figure 8: Full reasoning intensiveness pairwise comparison prompt template used in our experiments.

SYSTEM MESSAGE

You are an expert judge analyzing why AI models fail on reasoning-intensive questions. Compare the correct and incorrect answers to determine if the failure was primarily due to insufficient reasoning ability or lack of domain knowledge.

Reply in this exact format:

```
###EXPLANATION: <detailed analysis of why the low-reasoning model
failed>
###RESULTS: REASONING/KNOWLEDGE/BOTH/UNCLEAR
```

USER MESSAGE

You will be shown a question from an academic dataset, along with

- (1) a *CORRECT* answer from a high-reasoning model and
- (2) an *INCORRECT* answer from a low-reasoning model.

Your task is to analyze *why* the low-reasoning model failed.

Consider whether the failure is primarily due to:

- *REASONING*: Insufficient logical thinking, problem-solving ability, or step-by-step analysis
- *KNOWLEDGE*: Lack of domain knowledge (missing facts, formulas, concepts, procedures)
- *BOTH*: Significant deficiencies in both reasoning and knowledge
- *UNCLEAR*: Cannot determine the primary cause of failure

QUESTION

```
Context: {{context}}
Question: {{question}}
CORRECT ANSWER (from {{high_model}}):
{{high_full_response}}
INCORRECT ANSWER
(from {{low_model}}):
{{low_full_response}}
```

Analyze why the low-reasoning model failed. Was it primarily due to insufficient reasoning ability or lack of knowledge?

Results We show that both protocols agree that filtered instances require significantly more reasoning efforts than non-filtered instances from SCIREAS, with (a) showing 71% agreement in accuracy by LLMs with 78% human annotation agreement and (b) showing 91% agreement by LLMs with 90% human agreement, where human annotations are made by authors on 80 sampled tests for each protocol.

D Frontier Model API Evaluation Configuration

For OpenAI and xAI provided reasoning models, we apply generic "low" and "high" reasoning effort parameters with respect to official documentation where specificity on token budget is not allowed; for other reasoning models that allows thinking budgets as input (e.g. Gemini and Anthropic), we adopt "low" as definition introduced by LiteLLM, which corresponds to 1024 budget, and remove the constraint to allow for as many thinking tokens as the model needed to unleash full potential as "high" reasoning effort, corresponding to the highest reasoning effort from OpenAI and xAI models. For all frontier reasoning models, if not restricted, we set temperature=1, borrowed from OpenAI forced setting, and top-p=0.95, borrowed from recommended setting by Anthropic, with max generation length of 64K, as we observe no models tend to output more than 20K tokens. We log API pricing at the time of writing in Table 8.

Model	Input Price (\$ per 1M tokens)	Output Price (\$ per 1M tokens)
OpenAI models		
GPT-4.1-2025-04-14	2.00	8.00
o3-mini-2025-01-31	1.10	4.40
03-2025-04-16	2.00	8.00
o4-mini-2025-04-16	1.10	4.40
GPT-5-2025-08-07	1.25	10.00
GPT-oss-120B (Together AI)	0.15	0.60
DeepSeek models		
DeepSeek-V3-0324	0.14	0.28
DeepSeek-R1-0120	0.55	2.19
DeepSeek-R1-0528	0.55	2.19
Alibaba Qwen models (Together AI)		
Qwen3-32B	0.40	1.20
Qwen3-235B-2507	0.65	3.00
Google models		
Gemini-2.5-Pro-Preview-05-06	1.25	10.00
Meta models (Together AI)		
Llama-4-Maverick-17B-128E-Instruct-FP8	0.27	0.85
Anthropic models		
Claude-Sonnet-4-20250514	3.00	15.00

Table 8: Pricing (\$ per 1M tokens) for input and output across different LLM providers at the time of writing, without any discounts.

E Training / Evaluation Details

E.1 DISTILLATION FROM REASONING LLMS

To obtain high-performing reasoning models for study, we employ a distillation method that fine-tunes smaller models using Supervised Fine-tuning (SFT) on the CoT trajectories generated by large reasoning models, as it is more effective than reinforcement learning (RL) with the small models alone (DeepSeek-AI et al., 2025). Specifically, we consider the standard SFT framework for language models where the objective is to train a model f_{θ} to approximate a distribution over output

⁷https://docs.litellm.ai/docs/providers/anthropic#usage—thinking-reasoning_content

⁸https://community.openai.com/t/o3-mini-unsupported-parameter-temperature/1140846/3

https://docs.anthropic.com/en/docs/build-with-claude/extended-thinking#feature-compatibility

sequences y conditioned on input x, based on a dataset $\mathcal{D} = \{(x^i, y^i)\}_{i=1}^N$. For recent reasoning LLMs such as DeepSeek-R1, the output y consists of two main parts: a reasoning trace r and the actual output a. In practice, the reasoning traces are enclosed by keywords <think> and </think>, indicating the start and the end of the reasoning process. The model is trained with the standard SFT objective: $\mathcal{L}(\theta) = -\Sigma_{(x,y)\in\mathcal{D}} \Sigma_{t=1}^{|y|} \log p_{\theta}(y_t|y_{< t},x)$, where y_t is the t-th token and $y_{< t}$ is its prefix.

E.2 EXTENDED SETUP

E.2.1 TRAINING SETTINGS

We filter out instances with a token length greater than 4096. The models are trained for 5 epochs with a cosine learning rate scheduler, a maximum learning rate of 1e-5, and 3% warmup steps.

E.2.2 EVALUATION SETUP

The reasoning models could produce excessively long outputs, and may be prone to self-repetition with greedy decoding (DeepSeek-AI, 2025). In this work, unless otherwise specified, we apply greedy decoding on non-CoT fine-tuned models and top-p=0.95, temperature=0.6 on reasoning models, with a maximum generation length of 64K. From our preliminary studies, we observe that the setup generally reflects the best performance for both settings, and the decoding setup matches the recommended setup from recent efforts in large reasoning models, such as Llama-Nemotron (Bercovich et al., 2025a). Notably, for Qwen (Yang et al., 2024) models and their variants, we apply YaRN context extension (Peng et al., 2023) as recommended by the official model card (Team, 2024).

E.3 MATH VS. NON-MATH

E.3.1 FILTERING HEURISTICS

We label instances as math-needed if they contain explicit numeric quantities that typically imply computation. Importantly, numbers that appear solely within unit expressions (e.g., "cm²") or chemical formulas (e.g., "H₂O" or "NaCl") are not treated as indicators of math-related reasoning.

Specifically, a question is marked *Has-Math* when it includes

- 1. a signed or unsigned integer/decimal (e.g. 3, -2.5, 60, 9.81),
- 2. **not** embedded inside a word (so digits in H2O, COVID-19, IL-2... are ignored), and
- 3. optionally followed—without intervening letters—by any one of the unit strings listed in Fig. 11.

E.3.2 COT IMPROVEMENTS ON MATH VS. NON-MATH

In response to the question in §3.2, we use categorize instances from SCIREAS into Has-Math and No-Math, resulting in 8,527 cases identified as Has-Math and 4,757 as No-Math. We compute the micro accuracy on frontier models and plot the performance gains by increasing the thinking budget from low reasoning effort to high reasoning effort in Figure 10.

E.3.3 EFFECTS ON REASONING-FINE-TUNED MODELS

As shown in Table 1, Qwen-STEM and Qwen-Math both exhibit significant improvement over the base model on SCIREAS and SCIREAS-PRO.

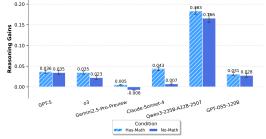


Figure 10: Performance gains on Has-Math instances vs. No-Math instances across different frontier models when reasoning effort increases. CI95% shown as the error bar. CoT helps more with Has-Math instances.

Qwen-Math slightly outperforms Qwen-STEM on SCIREAS and the gap is amplified on SCIREAS-PRO.

 $^{^{10}}$ Longer input lengths would slow down our training in quadratic order based on 8 80GB A100/H100 GPUs.

Units recognised by the heuristic								
 % °C, °F, K, ° g, kg, mg, μg/ug, lb/lbs, oz m, cm, mm, km L/l, mL/ml, μL/μl/ul Pa, kPa, MPa, atm, bar, mbar 	kW, MW, GW							

Figure 11: Unit suffixes accepted by the nu-
meric heuristic. A standalone number with any
of these units (or no unit) is treated as evidence
that the question contains mathematical content.

Model	Has-Math Acc.	No-Math Acc.							
SCIREAS-PRO: 1,260 Instances									
#	1,172	88							
Qwen	14.25	12.50							
Qwen-STEM	15.53	23.86							
Qwen-Math	17.58	13.64							
Qwen-BOTH	20.56	28.41							
Llama	11.52	13.64							
Llama-STEM	14.16	15.91							
Llama-Math	17.24	13.64							
Llama-BOTH	15.96	23.86							

Table 9: Accuracy breakdown on math and nonmath instances for SCIREAS-PRO. -Math and -STEM variants contribute to different dimensions of performance, while -BOTH captures improvements on both.

Given limited subject coverage on SYNTHETIC-1-Math dataset, the strong performance of check-points fine-tuned on it only seems surprising — Does the improvement come from generalization from math reasoning to a wider domain, or is it because the high-reasoning instances in our datasets are math-intensive? To answer this question, we categorize SCIREAS-PRO into math and non-math instances by heuristics.

As shown in Table 9, we find that math computation appears frequently among reasoning-intensive instances, and the improvements on SCIREAS-PRO mostly come from improved math capabilities. For non-math instances, -math variants hardly improve, while -STEM variants and -BOTH variants, trained with STEM subjects data, show noticeable improvements.

E.4 Training Knowledge Enhanced Scientific Reasoning Models

Our post-trained checkpoints are based on models fine-tuned on either SYNTHETIC-1-Math, SYNTHETIC-1-STEM, or both, while combining the two, which cover both STEM and mathematical reasoning, achieves the strongest performance (Table 1). To further assess the effectiveness of this Math+STEM data mixture following §4.1, we compare it directly against concurrently released long-CoT SFT datasets on the same base model. We then apply the same mixture to Qwen3-8B-Base to obtain SCILIT01 to provide a stronger baseline.

Specifically, we compare Qwen-BOTH, which is fine-tuned using our training recipe, with SYNTHETIC-1-SFT Mattern et al. (2025), a model fine-tuned on SYNTHETIC-1 with additional coding and preference alignment data, and Qwen-Nemotron, a model we trained with the same settings and same amount of data (§4.1) sampled from science and math domains of Llama-Nemotron Bercovich et al. (2025b), a training data mixture for reasoning fine-tuning, all post-trained on Qwen2.5-7B-Instruct. The results in Table 10 show that our data composition yields a stronger baseline for scientific reasoning than concurrent data recipes on Qwen2.5-7B-Instruct (Table 10 center block), and Qwen-BOTH reaches comparable performance to models from concurrent efforts focusing on reasoning enhancement post-training recipes (Table 10 left-hand block, i.e., OpenR1 Face (2025), Llama-Nemotron Bercovich et al. (2025b), and General-Reasoner Ma et al. (2025b)).

Furthermore, using our recipe, we fine-tune the recently released Qwen3-8B-Base to deliver a stronger model, SCILIT01. While its performance falls behind Qwen3-8B with the thinking mode, which has undergone more sophisticated post-training, it outperforms Qwen3-8B with non-thinking mode (Table 10 right-hand block). This indicates that SCILIT01 partially unleashes the reasoning capabilities from the base model, offering a strong baseline for future study on post-training recipe for scientific reasoning.

Models	OpenR1	Llama-Nemotron	General-Reasoner	SYNTHETIC-1-SFT	Qwen-Nemotron	Qwen-BOTH	SciLit01	Qwen3	Qwen3-thinking
Base Model	Q2.5-Math	L3.1-Inst.	Q2.5-Base	•	Q2.5-Inst			Q3-Base	
Training Methods	SFT	SFT&RL	RL	SFT	SFT	SFT	SFT	–	-
Trained by Us	No	No	No	No	Yes	Yes	Yes	No	No
GPQA SuperGPQA* MMLU-Pro* LabBench* OlympiadBench SciBench SciEval* SciKnowEval* SciRIFF*	44.42	37.95	35.94	38.84	44.20	40.63	50.89	55.80	55.80
	31.90	29.39	14.26	22.39	19.47	20.33	30.11	23.32	38.27
	60.86	65.64	62.14	56.21	63.57	65.00	76.57	73.36	81.71
	27.14	27.78	35.58	28.61	35.76	33.00	35.07	36.99	38.19
	53.03	37.62	19.82	40.75	29.33	34.55	43.78	28.51	21.30
	61.85	57.66	19.08	51.59	48.27	47.11	61.27	54.05	68.21
	43.64	68.67	70.34	46.41	38.53	72.36	80.60	81.51	84.02
	28.45	30.69	34.19	19.13	31.85	32.00	39.46	37.99	41.81
	29.17	34.01	37.75	28.57	39.24	41.81	44.01	47.23	47.26
UGPhysics* Average SCIREAS-PRO	50.30	45.92	20.86	43.96	46.52	40.03	52.28	30.98	59.81
	43.08	43.53	34.99	37.64	39.67	42.68	<u>51.41</u>	46.97	53.64
	26.43	23.75	13.73	19.44	19.68	21.11	24.84	19.05	29.92

Table 10: Performance of concurrent efforts on open-recipe post-training in <10B-parameter level. SCILIT01 shows competitive performance relative to concurrent reasoning post-training methods. We abbreviate Qwen2.5, Qwen3, and Llama-3.1 as Q2.5, Q3, and L3.1, respectively; '-Inst.' denotes the instruction-tuned variant. The best and second-best overall results are highlighted in bold and underlined, respectively.

F EXTENDED KRUX DETAILS

F.1 KNOWLEDGE EXTRACTION

In this work, we apply DeepSeek-R1 as the extractor. Prompt shown in Figure 12. We show a set of KIs extracted from Qwen2.5-7B-Instruct (Figure 13), Qwen-Math variants (Figure 14), and DeepSeek-R1 (Figure 15) for the same question from GPQA:

Question: A large gene has dozens of exons, of which the central ones code for folded triple helical repeats that connect the cytoskeleton with sarcolemma and extracellular space. Each exon usually codes for one folded triple alpha helix. The most common mutations of the gene are central exon deletions that create out-of-frame peptides and progressive degenerative organ waste. A solution is to deliver a Morpholino that recognizes the 5' end of the out-of-frame exon in pre-mRNA. The molecule prevents binding of the spliceosome and creates exon skipping and in-frame joining. Several missing exons are well tolerated by an organism. Which structure below is not involved in the proposed therapy? (A) lariat (B) antisense (C) R-loops (D) polyA tail.

1340

1341

1347 1348

USER MESSAGE

You are given a reasoning chain that explains and justifies a particular conclusion or answer. Your task is to extract **all distinct knowledge pieces** from this chain. A knowledge piece is any standalone statement conveying an explicit fact, definition, mechanism, relationship, or insight that can be generalized beyond the specific question.

Instructions:

- 1. Read the entire reasoning chain.
- 2. Identify each discrete fact or insight expressed.
- 3. Rewrite each as a self-contained, generalizable sentence.
- 4. Do **not** include any contextual or example-specific details.
- 5. Output **only** a list of those sentences.

Output Format:

- knowledge-Piece-1
- knowledge-Piece-2
- .

Reasoning Chain: {{REASONING}}

Now perform the extraction.

Figure 12: Full knowledge extraction prompt template used in our experiments.

KNOWLEDGE EXTRACTION OUTPUT

Based on the reasoning chain, here are all distinct knowledge pieces:

- R-loops are structures formed when a strand of RNA hybridizes with its complementary DNA strand, displacing the other DNA strand.
- 2. Morpholino can prevent spliceosome binding.
- 3. Preventing spliceosome binding can lead to exon skipping.
- 4. Exon skipping can result in in-frame joining of exons.
- 5. The process of using Morpholino to induce exon skipping does not involve R-loops.

Figure 13: Example KIs extracted from Qwen2.5-7B-Instruct.

KNOWLEDGE EXTRACTION OUTPUT

Based on the reasoning chain, here are all distinct knowledge pieces:

- Morpholino oligonucleotides can bind to the 5' splice site of pre-mRNA to block spliceosome recognition.
- 2. Blocking spliceosome binding prevents normal intron removal and leads to exon skipping.
- 3. Exon skipping can enable in-frame joining of adjacent exons, preserving the reading frame.
- Antisense molecules function by complementary base-pairing to specific RNA sequences to modulate splicing.
- Lariat structures are formed during canonical splicing via branch site interactions but are absent when spliceosome activity is inhibited.
- R-loops are DNA-RNA hybrid structures unrelated to spliceosome-mediated exon skipping mechanisms.
- The polyA tail is added post-transcriptionally and is not directly involved in splice site selection or exon skipping processes.

Figure 14: Example KIs extracted from Qwen-Math.

KNOWLEDGE EXTRACTION OUTPUT

1350

1351 1352

1353

1354

1355

1356

1357

1358

1359

1363

1365

1367

1369

1370

1371

1372

1373

1374

1375

1376

1379 1380

1381

1382

1386

1387

1388

1389 1390

1391 1392

1393

1394

1400 1401

1402

1403

Based on the reasoning chain, here are all distinct knowledge pieces:

- 1. Morpholinos are antisense oligonucleotides that bind to complementary RNA sequences.
- The spliceosome recognizes splice sites on pre-mRNA through components like U1 snRNP, which binds to the 5' splice site.
- Blocking spliceosome access to a splice site prevents recognition of exon boundaries and causes exon skipping.
- 4. Exon skipping can restore the reading frame by joining exons flanking the skipped exon.
- Functional proteins may be produced if skipped exons encode non-critical regions and the remaining exons maintain an in-frame sequence.
- 6. Lariat structures form during canonical splicing when introns are removed via branch site cleavage.
- R-loops are nucleic acid structures involving RNA-DNA hybrids and are not directly involved in spliceosome-mediated splicing.
- 8. Antisense oligonucleotide therapies rely on sequence-specific binding to pre-mRNA targets.
- Polyadenylation (polyA tail addition) occurs during mRNA maturation and is unrelated to splice site selection or exon skipping mechanisms.

Figure 15: Example KIs extracted from DeepSeek-R1.

F.2 ALTERNATIVE EXTRACTOR

We generate KIs from the same set of knowledge traces as the knowledge source and inspect them in pairs blindly. Specifically, we compare Qwen3-30B-A3B-Thinking-2507 and DeepSeek-R1 as extractors. Among 100 sample pairs, *all* KIs generated from both models (a) do not refer to any specific identities in the question, e.g., "The option B ...", (b) do not expose the final answers, and (c) adhere to the given traces with no additional information nor missing essential components. However, different extractors may provide KIs at different granularities. Specifically, we observe that the DeepSeek-R1 extractor tends to be more detailed in terms of KIs extracted, as shown in the example in Figure 16. Our human annotation shows that DeepSeek-R1 is more detailed, with 53% of sampled pairs, 43% tied, and only 4% of Qwen3-30B-A3B being more detailed.

KI EXAMPLE FROM DeepSeek-R1

- Copulation refers to the physical act of mating between male and female organisms.
- Insemination is the process of depositing semen containing sperm into the female reproductive tract.
- Fertilization is the union of a sperm cell and an egg cell to form a zygote.
- Copulation encompasses the physical joining of organisms, while insemination specifically involves sperm transfer.
- Fertilization is a distinct biological event from copulation and insemination, occurring after successful insemination.
- Copulation and fertilization are separate processes, with the former involving physical mating and the latter involving gamete fusion.

KI EXAMPLE FROM Qwen3-30B-A3B-Thinking-2507

- Copulation is the general term for the act of mating.
- Copulation is the physical joining of the male and female.
- Insemination is the male depositing sperm into the female's reproductive tract.
- Fertilization is the union of gametes.

Figure 16: An example of KIs generated by DeepSeek-R1 and Qwen3-30B-A3B-Thinking-2507. The DeepSeek-R1 tends to cover knowledge recalled in more detail.

F.3 KNOWLEDGE PROBING

We provide our probing question synthesis prompt (Figure 17), example input and output (Figure 18), and knowledge probing results in Table 4.

USER MESSAGE You are a meticulous question-authoring assistant. Your job is to convert declarative knowledge statements into *probing* multiple-choice questions that can test whether another language model truly stores the fact in its parametric memory. ## IMPORTANT INSTRUCTIONS FOR QUESTIONS: 1. Factual: It should be about a specific detail or fact mentioned in the statement. For example, a true or false statement, a statistic, a definition, etc. 2. Important: It should be a question about the main topic or a key detail/finding/conclusion of the statement. 3. Context-Independent: It should be fully understandable on its own, without phrases like "the pro-posed model" or "this approach" that assume prior context. ## IMPORTANT INSTRUCTIONS FOR ANSWERS: 1. Provide one correct answer and 4 - 6 incorrect answers. 2. Ensure all answers are roughly the same length and follow a similar style so the correct answer cannot be guessed based on length or style alone. 3. The incorrect answers must be plausible but ultimately wrong, reflecting a misunderstanding or misinterpretation of the knowledge. ## OUTPUT FORMAT: Please provide the question, correct answer, incorrect answers, and a list of text snippets from the article as "evidences" in the following format: { "question": "Your question here", "correct_answer": "Correct answer here", "incorrect_answers": ["Incorrect answer 1", ..., "Incorrect answer N"], "evidences": ["Text snippets from the article that supports the question and correct answer", "Another text snippet"]

Figure 17: Knowledge probing question synthesis template used in our experiments.

Please provide your response in the specified format without any additional text.

EXAMPLE src_text

Knowledge Statement: {src_text}

"Hyperfine structure in EPR spectroscopy arises from interactions between unpaired electrons and nuclear spins."

EXAMPLE OUTPUT

```
{
  "question": "What causes hyperfine structure in EPR spectroscopy?",
  "correct_answer": "Interactions between unpaired electrons and nuclear spins",
  "incorrect_answers": [
  "Interactions between electron spins and lattice vibrations", "Coupling between electron orbitals and magnetic fields", "Dipolar interactions between neighboring nuclei", "Spin-orbit coupling within the electron cloud", "Chemical shift anisotropy of atomic orbitals"],
  "evidences": [
  "Hyperfine structure in EPR spectroscopy arises from interactions between unpaired electrons and nuclear spins."]
}
```

Figure 18: Knowledge probing question synthesis example input and output.

G LLM USAGE STATEMENT

We used ChatGPT-o3 from OpenAI for grammar and typo corrections.

Domain	Task Source	Subtask/Subdomain	Instances	Total	Metrics
	GPQA	Physics	187		Acc
	MMLU-Pro	physics	200		Acc
	SciBench	fund	81		Acc
		thermo	83		Acc
		class	63		Acc
	OlympiadBench-COMP	physics_en	236		Acc
	SciKnowEval.L5	physics_problem_solving	200		LM
	SciEval	physics_knowledge_application	29		Acc
		physics_scientific_calculation	200		Acc
	UGPhysics	Electrodynamics	170		Acc
	•	Thermodynamics	200		Acc
Physics		GeometricalOptics	54	5087	Acc
<i>J</i>		Relativity	200		Acc
		ClassicalElectromagnetism	200		Acc
		ClassicalMechanics	200		Acc
		WaveOptics	200		Acc
		QuantumMechanics	200		Acc
		TheoreticalMechanics	200		Acc
		AtomicPhysics	200		Acc
		SemiconductorPhysics	148		Acc
		Solid-StatePhysics	154		Acc
		StatisticalMechanics	200		Acc
	SuperGPQA	Physics	1482		Acc
	GPQA	Chemistry	183		Acc
	MMLU-Pro	chemistry	200		Acc
	SciBench	quan	41		Acc
Chemistry		chemc	47		Acc
		atkins	121		Acc
		matter	57	2158	Acc
	SciKnowEval.L5	chemical_procedure_generation	74		LM
		chemical_reagent_generation	125		LM
	SciEval	chemistry_knowledge_application	200		Acc
		chemistry_scientific_calculation	200		Acc
	SuperGPQA	Chemistry	910		Acc
Comp Sci	MMLU-Pro	computer science	200	415	Acc
	SciRIFF	Qasper	107	415	F1, LM
	SuperGPQA	Computer Science and Technology	108		Acc
Math	MMLU-Pro	math	200		Acc
	SciBench	calc	52		Acc
		stat	92	2522	Acc
		diff	55	2533	Acc
	OlympiadBench-COMP	maths_en	674		Acc
	SuperGPQA	Mathematics	1460		Acc

Table 11: Domain-wise breakdown of SCIREAS tasks and instance counts (Part 1: Physics to Math).

Domain	Task Source	Subtask	Instances	Total	Metrics
	GPQA	Biology	78		Acc
Biology	MMLU-Pro	biology	200		Acc
	LabBench	CloningScenarios	33		Acc
		ProtocolQA	108		Acc
		SeqQA	600	1911	Acc
	SciKnowEval.L5	biological_procedure_generation	200	1711	LM
	a	biological_reagent_generation	200		LM
	SciEval	biology_knowledge_application	200		Acc
	C CDO A	biology_scientific_calculation	200		Acc
	SuperGPQA	Biology	92		Acc
	MMLU-Pro	health	200		Acc
Medicine	SciRIFF	SciFact	184	634	F1, LM
		Evidence Inference	250		F1
	SciKnowEval.L5	crystal_structure_and_composition	196		LM
		specified_band_gap_material_generation	200	624	LM
Material Sci		property_and_usage_analysis	118		LM
	SuperGPQA	Materials Science and Engineering	110		Acc
	MMLU-Pro	engineering	200		Acc
	SuperGPQA	Control Science and Engineering	77		Acc
		Information and Communication En-	156		Acc
		gineering			
Engineering		Electrical Engineering	234		Acc
		Chemical Engineering and Technol-	226		Acc
		ogy Power Engineering and Engineering	345	2205	Acc
		Thermophysics	373		Acc
		Electronic Science and Technology	95		Acc
		Hydraulic Engineering	67		Acc
		Mechanics	456		Acc
		Mechanical Engineering	30		Acc
		Civil Engineering	93		Acc
		Optical Engineering	162		Acc
		Nuclear Science and Technology	30		Acc
		Instrument Science and Technology	12		Acc
		Systems Science	22		Acc

Table 12: Domain-wise breakdown of SCIREAS tasks and instance counts (Part 2: Biology to Engineering).

Domain	Task Source	Subtask/Subdomain	Instances	Total	Metrics
	GPQA	Physics	8		Acc
	MMLU-Pro	physics	5		Acc
	SciBench	fund	1		Acc
		thermo	10		Acc
		class	8		Acc
	OlympiadBench-COMP	physics_en	25		Acc
	SciEval	physics_knowledge_application	1		Acc
		physics_scientific_calculation	1		Acc
	UGPhysics	Electrodynamics	17		Acc
		Thermodynamics	16		Acc
Physics		GeometricalOptics	9	388	Acc
1 Hysics		Relativity	16	300	Acc
		ClassicalElectromagnetism	21		Acc
		ClassicalMechanics	17		Acc
		WaveOptics	16		Acc
		QuantumMechanics	17		Acc
		TheoreticalMechanics	13		Acc
		AtomicPhysics	13		Acc
		SemiconductorPhysics	13		Acc
		Solid-StatePhysics	13		Acc
		StatisticalMechanics	15		Acc
	SuperGPQA	Physics	133		Acc
	GPQA	Chemistry	31		Acc
	MMLU-Pro	chemistry	3		Acc
Chemistry	SciBench	quan	3		Acc
		chemc	2		Acc
		atkins	6	135	
		matter	3		Acc
	SciEval	chemistry_knowledge_application	11		Acc
		chemistry_scientific_calculation	3		Acc
	SuperGPQA	Chemistry	73		Acc
Comp Sci	MMLU-Pro	computer science	6	21	Acc
1	SuperGPQA	Computer Science and Technology	15		Acc
Math	MMLU-Pro	math	3		Acc
	SciBench	calc	2		Acc
		stat	2	202	Acc
		diff	3	283	Acc
	OlympiadBench-COMP	maths_en	92		Acc
	SuperGPQA	Mathematics	181		Acc

Table 13: Domain-wise breakdown of SCIREAS-PRO tasks and instance counts (Part 1: Physics to Math).

Domain	Task Source	Subtask	Instances	Total	Metrics
	GPQA	Biology	2		Acc
	MMLU-Pro	biology	6		Acc
	LabBench	CloningScenarios	2		Acc
Biology		ProtocolQA	10	123	Acc
		SeqQA	89	123	Acc
	SciEval	biology_knowledge_application	3		Acc
	a ano.	biology_scientific_calculation	2		Acc
	SuperGPQA	Biology	9		Acc
Medicine	MMLU-Pro	health	5	5	Acc
Material Sci	SuperGPQA	Materials Science and Engineering	13	13	Acc
	MMLU-Pro	engineering	14		Acc
	SuperGPQA	Control Science and Engineering	7		Acc
		Information and Communication En-	15		Acc
		gineering			
		Electrical Engineering	32		Acc
		Chemical Engineering and Technol-	43		Acc
		ogy		202	
Engineering		Power Engineering and Engineering Thermophysics	44	292	Acc
		Electronic Science and Technology	13		Acc
		Hydraulic Engineering	13		Acc
		Mechanics	54		Acc
		Mechanical Engineering	7		Acc
		Civil Engineering	18		Acc
		Optical Engineering	23		Acc
		Nuclear Science and Technology	3		Acc
		Instrument Science and Technology	2		Acc
		Systems Science	4		Acc

Table 14: Domain-wise breakdown of SCIREAS-PRO tasks and instance counts (Part 2: Biology to Engineering).