

# DEMYSTIFYING SCIENTIFIC PROBLEM-SOLVING IN LLMs BY PROBING KNOWLEDGE AND REASONING

**Anonymous authors**

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.<sup>1</sup>

## 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,

<sup>1</sup>The codebase and artifacts are released at [link-redacted-for-review](#).

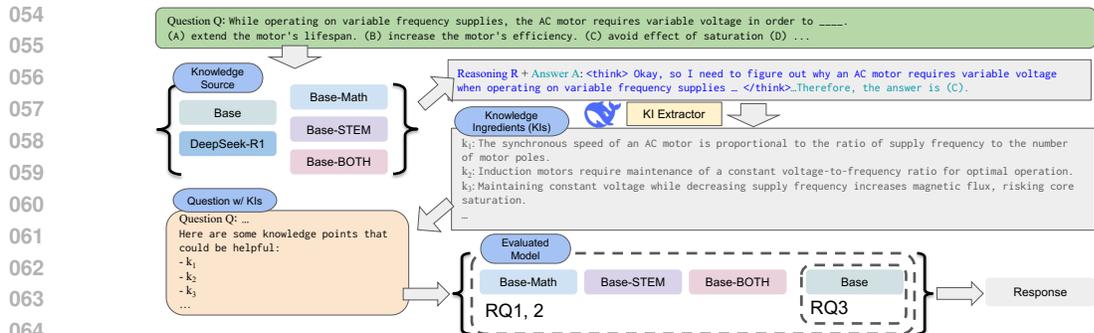


Figure 1: KRUX pipeline. Starting upper left, we prompt an LLM (one of base, DeepSeek-R1, Base-Math, Base-STEM, and Base-BOTH) with a question from SCIREAS as the knowledge source, collect the output and reasoning traces, and feed the reasoning traces to DeepSeek-R1 as the extractor to generate knowledge 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** (**K**nowledge & **R**easoning **U**tization **e**Xams), 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.

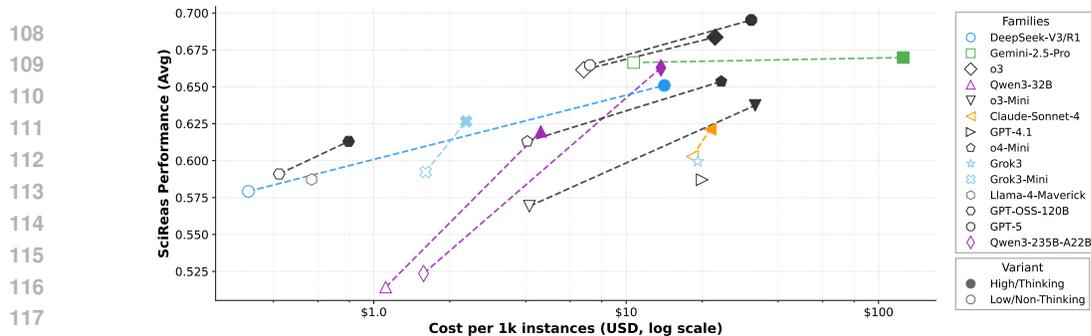


Figure 2: Frontier reasoning models’ performance evaluated on SciREAS. The X-axis shows the cost per 1k instances in USD. Different reasoning settings on the same model can result in distinct costs and performance, but the margins vary depending on the models.

## 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 multi-step reasoning or standardized implementation. For example, most tasks in SciRIF (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. Benchmarks like CURIE (Kon et al., 2025) and QASA (Lee et al., 2023) offer exposure to scientific subjects, but their long-context feature hinders our objective to disentangle the effect of knowledge and reasoning. 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), SciRIF, UGPhysics (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.

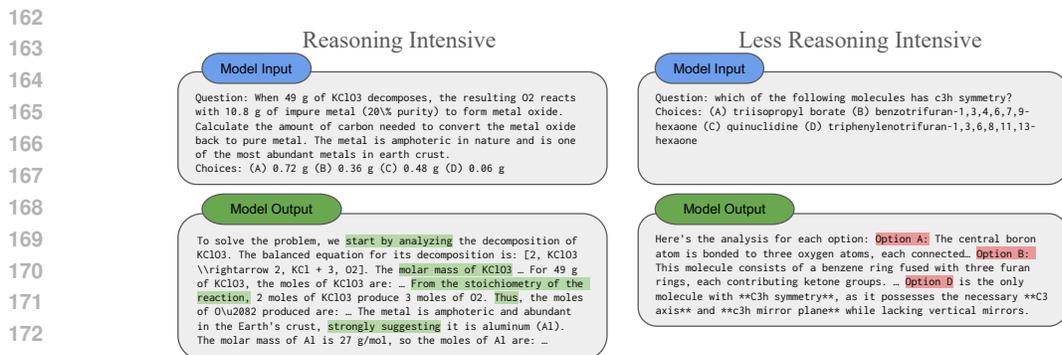


Figure 3: An example pair with varying reasoning intensity, where the example on the left is sampled from SCIREAS-PRO and the right is a filtered out example (§3.1). On the left, the progressive reasoning chain is highlighted. The example on the right emphasizes knowledge recall on each option with a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

### 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 a unified implementation optimized with vLLM (Kwon et al., 2023) and batch job APIs<sup>2</sup> 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 (1) select candidate tasks based on documentation in their associated manuscripts that describe each subtask’s characteristics, which indicate its difficulty and reasoning intensity. After that, we (2) use a subtask-level exclusion protocol to retain only those tasks that require both domain knowledge and multi-step reasoning.<sup>3</sup> Our procedure involves 3 authors as annotators, where two authors decide jointly in the first round, and another author validates by conducting the filtering process independently, following the same selection policy. This two-fold annotation validation reached an agreement accuracy of 90.1%. We provide more details and explanations for our protocol and the subtasks we selected 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

<sup>2</sup>We 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.

<sup>3</sup>While this manual inspection can be subjective, it is based on the authors’ graduate-level expertise. We provide more details and examples in Appendix C.1.

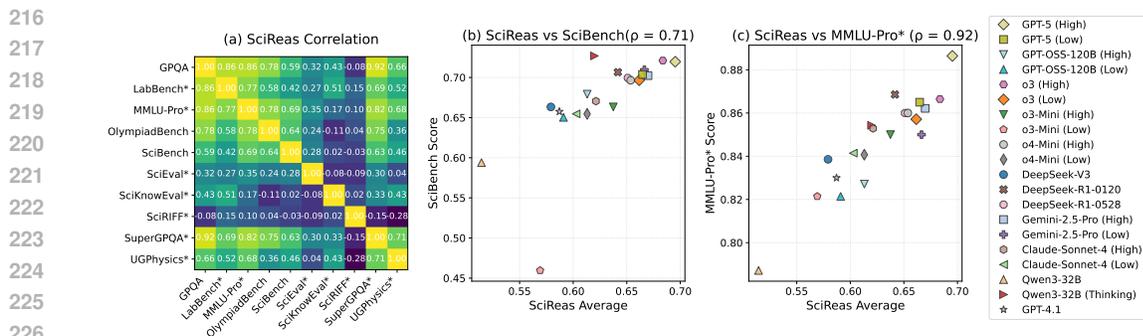


Figure 4: SCIREAS correlations breakdown. (a) Task-to-task Pearson correlations. SCIREAS incorporates tasks complementary to popular benchmarks. (b) and (c) show performance on SCIREAS vs. SciBench and MMLU-Pro\*. Models may be tuned for certain tasks, outperforming higher-ranked models on individual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

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).<sup>4</sup> 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.

### 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.<sup>5</sup>

**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.

<sup>4</sup>Because 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).

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

**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 ( $>10\times$ ) 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 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 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.

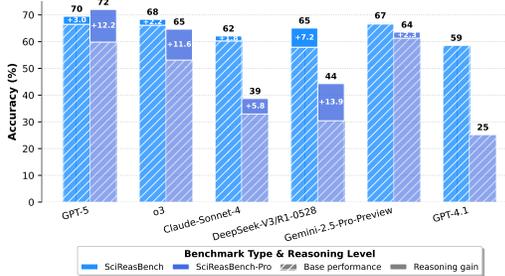


Figure 5: Model performance on SCIREAS and SCIREAS-PRO with varying reasoning capabilities. SCIREAS-PRO amplifies gaps between low-reasoning and high-reasoning settings.

#### 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-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

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	SCI REAS	-PRO
<i>Our Checkpoints</i>			
Qwen	–	37.07	13.97
Qwen-STEM	SFT	40.47	16.11
Qwen-Math	SFT	41.99	18.17
Qwen-BOTH	SFT	<b>42.84</b>	<b>21.11</b>
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
<i>Concurrent Reasoning Post-training</i>			
SYNTHETIC-1-SFT	SFT	37.64	19.44
OpenR1	SFT	43.08	<b>26.43</b>
Llama-Nemotron	SFT&RL	<b>43.53</b>	23.75
General-Reasoner	RL	34.99	13.73

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.<sup>6</sup> 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 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). A KI is any standalone statement conveying an explicit fact, definition, mechanism, relationship, or insight that can be generalized beyond the specific question. It is a self-contained, generalizable sentence and does not include any contextual or example-specific details. The final prompt we use is carefully tuned based on our manual inspections of the extracted KIs (provided in Appendix F.1 with examples). We then augment the original question by prepending the extracted set of KIs in-context and ask the models to solve the same problem.

To validate the extraction quality and generalizability, we perform additional checks on DeepSeek-R1 and Qwen3-30B-A3B-Thinking-2507 as the extractors to ensure that the 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 to in option B ...”), (b) do *not* leak any part of the final answer, and (c) strictly adhere to the traces as the knowledge source without adding extra 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, and provide more details and experiment results with alternative extractors 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/

<sup>6</sup>Notably, 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.

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:** To what extent 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 fine-tuning. The fact that a base model without strong reasoning capabilities can outperform reasoning models in this setting suggests that their parametric knowledge lacks the information in the KIs, or that they struggle to retrieve it from their parametric storage in a way that hinders 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	$34.24 \pm 0.93$	$30.93 \pm 1.43$
w/ R1 KIs	<b><u><math>47.19 \pm 1.53</math></u></b>	<b><u><math>41.40 \pm 2.46</math></u></b>
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	$29.06 \pm 1.44$	$34.40 \pm 2.58$
w/ R1 KIs	<b><u><math>43.57 \pm 0.88</math></u></b>	<b><u><math>42.27 \pm 1.60</math></u></b>
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: TO WHAT EXTENT DO REASONING-ENHANCED MODELS BENEFIT FROM EXTERNAL KNOWLEDGE?

**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.

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.

Models	GPQA		MMLU-Pro*		LabBench*	
	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	<b>41.63 ± 2.10</b>	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	<b>66.93 ± 0.72</b>	<b>74.00 ± 0.59</b>	30.18 ± 1.65	41.17 ± 2.32
Qwen-BOTH	40.81 ± 2.04	<b>54.46 ± 1.27</b>	65.71 ± 0.74	71.64 ± 1.16	<b>33.83 ± 2.59</b>	<b>43.90 ± 2.71</b>
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	<b>39.43 ± 2.00</b>	<b>54.73 ± 1.75</b>	<b>63.81 ± 0.90</b>	<b>72.74 ± 0.26</b>	<b>36.61 ± 2.73</b>	<b>48.65 ± 0.49</b>

**Answer to RQ2: Reasoning models also substantially benefit from the 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 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.

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 <i>Qwen-Math</i>	-Math <i>Llama-Math</i>	Llama <i>Llama-Math</i>	-Math
KI-GPQA	72.30	73.02	70.94	68.94
KI-MMLU-Pro*	82.49	81.50	74.46	74.12

**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	34.24 ± 0.93	59.03 ± 0.34
	w/ Qwen-Math KIs	<b>36.93 ± 1.75</b>	<b>63.66 ± 0.45</b>
Llama	w/ Llama KIs	29.06 ± 1.44	47.73 ± 0.89
	w/ Llama-Math KIs	<b>29.69 ± 1.72</b>	<b>53.91 ± 0.94</b>

## 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

- Iván Arcuschin, Jett Janiak, Robert Krzyzanowski, Senthoooran 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 Reakesh, 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 Norrick, 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,

- 540 Sandip Bhaskar, Maka Dong, Nave Assaf, Shahar Mor, Omer Ullman Argov, Scot Junkin, Olek-  
541 sandr Romanenko, Pedro Larroy, Monika Katariya, Marco Rovinelli, Viji Balas, Nicholas Edel-  
542 man, Anahita Bhiwandiwalla, Muthu Subramaniam, Smita Ithape, Karthik Ramamoorthy, Yut-  
543 ing Wu, Suguna Varshini Velury, Omri Almog, Joyjit Daw, Denys Fridman, Erick Galinkin,  
544 Michael Evans, Shaona Ghosh, Katherine Luna, Leon Derczynski, Nikki Pope, Eileen Long,  
545 Seth Schneider, Guillermo Siman, Tomasz Grzegorzec, Pablo Ribalta, Monika Katariya, Chris  
546 Alexiuk, Joey Conway, Trisha Saar, Ann Guan, Krzysztof Pawelec, Shyamala Prayaga, Oleksii  
547 Kuchaiev, Boris Ginsburg, Oluwatobi Olabiyi, Kari Briski, Jonathan Cohen, Bryan Catanzaro,  
548 Jonah Alben, Yonatan Geifman, and Eric Chung. Llama-nemotron: Efficient reasoning models,  
549 2025a. URL <https://arxiv.org/abs/2505.00949>.
- 550 Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil,  
551 Zach Moshe, Tomer Ronen, Najeeb Nabwani, et al. Llama-nemotron: Efficient reasoning models.  
552 *arXiv preprint arXiv:2505.00949*, 2025b.
- 553 Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. Crawling the internal knowledge-base  
554 of language models, 2023. URL <https://arxiv.org/abs/2301.12810>.
- 555 DeepSeek-AI. Deepseek-r1: Usage recommendations, 2025. URL <https://huggingface.co/deepseek-ai/DeepSeek-R1#usage-recommendations>.
- 556 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,  
557 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,  
558 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao  
559 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,  
560 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,  
561 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,  
562 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang  
563 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai  
564 Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,  
565 Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang,  
566 Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang,  
567 Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang,  
568 R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng  
569 Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing  
570 Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen  
571 Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong  
572 Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu,  
573 Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xi-  
574 aosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia  
575 Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng  
576 Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong  
577 Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong,  
578 Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou,  
579 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying  
580 Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda  
581 Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu,  
582 Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu  
583 Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforce-  
584 ment learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 585 Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL <https://github.com/huggingface/open-r1>.
- 586 Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua  
587 Yao, Qiang Zhang, and Huajun Chen. Sciknoweval: Evaluating multi-level scientific knowledge  
588 of large language models, 2024. URL <https://arxiv.org/abs/2406.09098>.
- 589 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,  
590 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-  
591 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang

- 594 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
595 evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.  
596
- 597 Daniela Gottesman and Mor Geva. Estimating knowledge in large language models without gener-  
598 ating a single token, 2024. URL <https://arxiv.org/abs/2406.12673>.  
599
- 600 Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom  
601 Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici,  
602 Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat,  
603 Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan,  
604 Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary  
605 Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an  
606 ai co-scientist, 2025. URL <https://arxiv.org/abs/2502.18864>.  
607
- 608 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
609 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd  
610 of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 611 Yuling Gu, Oyvind Tafjord, Bailey Kuehl, Dany Haddad, Jesse Dodge, and Hannaneh Hajishirzi.  
612 Olmes: A standard for language model evaluations, 2025. URL <https://arxiv.org/abs/2406.08446>.  
613
- 614 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han,  
615 Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. OlympiadBench:  
616 A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scien-  
617 tific problems. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the*  
618 *62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-*  
619 *pers)*, pp. 3828–3850, Bangkok, Thailand, August 2024. Association for Computational Linguis-  
620 tics. doi: 10.18653/v1/2024.acl-long.211. URL <https://aclanthology.org/2024.acl-long.211/>.  
621
- 622 Mingyu Jin, Weidi Luo, Sitao Cheng, Xinyi Wang, Wenyue Hua, Ruixiang Tang, William Yang  
623 Wang, and Yongfeng Zhang. Disentangling memory and reasoning ability in large language  
624 models, 2025. URL <https://arxiv.org/abs/2411.13504>.
- 625 Jude Khouja, Karolina Korgul, Simi Hellsten, Lingyi Yang, Vlad Neacsu, Harry Mayne, Ryan  
626 Kearns, Andrew Bean, and Adam Mahdi. Lingoly-too: Disentangling reasoning from knowl-  
627 edge with templatised orthographic obfuscation, 2025. URL <https://arxiv.org/abs/2503.02972>.  
628
- 629 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
630 language models are zero-shot reasoners, 2023. URL <https://arxiv.org/abs/2205.11916>.  
631
- 632 Patrick Tser Jern Kon, Jiachen Liu, Qiuyi Ding, Yiming Qiu, Zhenning Yang, Yibo Huang, Jayanth  
633 Srinivasa, Myungjin Lee, Mosharaf Chowdhury, and Ang Chen. Curie: Toward rigorous and  
634 automated scientific experimentation with ai agents, 2025. URL <https://arxiv.org/abs/2502.16069>.  
635
- 636 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
637 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
638 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*  
639 *Systems Principles*, 2023.  
640
- 641 Nathan Lambert, Lewis Tunstall, Nazneen Rajani, and Tristan Thrush. Huggingface h4  
642 stack exchange preference dataset, 2023. URL <https://huggingface.co/datasets/HuggingFaceH4/stack-exchange-preferences>.  
643
- 644 Jon M Laurent, Joseph D Janizek, Michael Ruzo, Michaela M Hinks, Michael J Hammer-  
645 ling, Siddharth Narayanan, Manvitha Ponnampati, Andrew D White, and Samuel G Rodrigues.  
646 Lab-bench: Measuring capabilities of language models for biology research. *arXiv preprint*  
647 *arXiv:2407.10362*, 2024.

- 648 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Chan Ho So, and Jaewoo Kang.  
649 Biobert: a pre-trained biomedical language representation model for biomedical text mining.  
650 *Bioinformatics*, 36(4):1234–1240, 2020.
- 651
- 652 Yoonjoo Lee, Kyungjae Lee, Sunghyun Park, Dasol Hwang, Jaehyeon Kim, Hong-in Lee, and  
653 Moontae Lee. Qasa: Advanced question answering on scientific articles. In *Proceedings of*  
654 *the 40th International Conference on Machine Learning*, 2023.
- 655 Aochong Oliver Li and Tanya Goyal. Memorization vs. reasoning: Updating llms with new knowl-  
656 edge, 2025. URL <https://arxiv.org/abs/2504.12523>.
- 657
- 658 Sihang Li, Jin Huang, Jiayi Zhuang, Yaorui Shi, Xiaochen Cai, Mingjun Xu, Xiang Wang, Lin-  
659 feng Zhang, Guolin Ke, and Hengxing Cai. Scilitlm: How to adapt llms for scientific literature  
660 understanding, 2025. URL <https://arxiv.org/abs/2408.15545>.
- 661 Qianchu Liu, Sheng Zhang, Guanghui Qin, Timothy Ossowski, Yu Gu, Ying Jin, Sid Kiblawi, Sam  
662 Preston, Mu Wei, Paul Vozila, Tristan Naumann, and Hoifung Poon. X-reasoner: Towards gener-  
663 alizable reasoning across modalities and domains, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2505.03981)  
664 [2505.03981](https://arxiv.org/abs/2505.03981).
- 665 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:  
666 Towards fully automated open-ended scientific discovery, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2408.06292)  
667 [abs/2408.06292](https://arxiv.org/abs/2408.06292).
- 668
- 669 Xinxu Lyu, Michael Duan, Rulin Shao, Pang Wei Koh, and Sewon Min. Frustratingly simple retrieval  
670 improves challenging, reasoning-intensive benchmarks, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2507.01297)  
671 [abs/2507.01297](https://arxiv.org/abs/2507.01297).
- 672 Ruotian Ma, Peisong Wang, Cheng Liu, Xingyan Liu, Jiaqi Chen, Bang Zhang, Xin Zhou, Nan Du,  
673 and Jia Li. S<sup>2</sup>r: Teaching llms to self-verify and self-correct via reinforcement learning, 2025a.  
674 URL <https://arxiv.org/abs/2502.12853>.
- 675
- 676 Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhui Chen. General-  
677 reasoner: Advancing llm reasoning across all domains. 2025b. URL [https://api.](https://api.semanticscholar.org/CorpusID:278768680)  
678 [semanticscholar.org/CorpusID:278768680](https://api.semanticscholar.org/CorpusID:278768680).
- 679 Justus Mattern, Felix Gabriel, and Johannes Hagemann. Synthetic-1 release: Two million collab-  
680 oratively generated reasoning traces from deepseek-r1, February 2025. URL [https://www.](https://www.primeintellect.ai/blog/synthetic-1-release)  
681 [primeintellect.ai/blog/synthetic-1-release](https://www.primeintellect.ai/blog/synthetic-1-release).
- 682
- 683 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke  
684 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time  
685 scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- 686 OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden  
687 Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko,  
688 Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally  
689 Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich,  
690 Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghor-  
691 bani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao,  
692 Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary  
693 Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang,  
694 Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel  
695 Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson,  
696 Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Eliz-  
697 abeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang,  
698 Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred  
699 von Lohman, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace  
700 Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart An-  
701 drin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichen,  
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,

- 702 Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish,  
703 Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan  
704 Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl  
705 Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu,  
706 Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam  
707 Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kon-  
708 draciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen,  
709 Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufner, Max Schwarzer, Meghan Shah, Mehmet  
710 Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael  
711 Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles  
712 Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil  
713 Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg  
714 Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov,  
715 Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar  
716 Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan  
717 Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agar-  
718 wal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu,  
719 Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph  
720 Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Tay-  
721 lor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson,  
722 Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna  
723 Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi  
724 Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen,  
725 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>.
- 726 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window  
727 extension of large language models, 2023. URL <https://arxiv.org/abs/2309.00071>.
- 728 Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao,  
729 James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim  
730 Rocktäschel, and Sebastian Riedel. Kilt: a benchmark for knowledge intensive language tasks,  
731 2021. URL <https://arxiv.org/abs/2009.02252>.
- 732 Arvind Prabhakar et al. Omniscience: A domain-specialized llm for scientific reasoning. *arXiv*  
733 *preprint arXiv:2503.17604*, 2025.
- 734 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien  
735 Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a  
736 benchmark. In *First Conference on Language Modeling (COLM)*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- 737 Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu,  
738 Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants,  
739 2025. URL <https://arxiv.org/abs/2501.04227>.
- 740 Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann  
741 Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To cot or not to cot? chain-  
742 of-thought helps mainly on math and symbolic reasoning. *ArXiv*, abs/2409.12183, 2024. URL  
743 <https://api.semanticscholar.org/CorpusID:272708032>.
- 744 Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu.  
745 Scieval: A multi-level large language model evaluation benchmark for scientific research, 2024.  
746 URL <https://arxiv.org/abs/2308.13149>.
- 747 Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun  
748 Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with  
749 llms. *arXiv preprint arXiv:2501.12599*, 2025a.
- 750 P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu,  
751 Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixin Deng, Shawn Gavin, Shian  
752 Jia, Sichao Jiang, Yiyao Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, David Ma,  
753

- 756 Yuansheng Ni, Haoran Que, Qiyao Wang, Zhoufutu Wen, Siwei Wu, Tyshawn Hsing, Ming  
757 Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin  
758 Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang,  
759 Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin,  
760 Yinghao Ma, Tianyang Pang, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu,  
761 Shanghaoran Quan, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo  
762 Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jinyang  
763 Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua  
764 Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Qian Liu, Tianyu Liu, Shiwen Ni, Junran Peng,  
765 Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang  
766 Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and  
767 Ge Zhang. Supergpqa: Scaling llm evaluation across 285 graduate disciplines, 2025b. URL  
768 <https://arxiv.org/abs/2502.14739>.
- 769 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.  
770 github.io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 771 Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL  
772 <https://qwenlm.github.io/blog/qwq-32b/>.
- 773
- 774 Rahul Thapa, Qingyang Wu, Kevin Wu, Harrison Zhang, Angela Zhang, Eric Wu, Haotian Ye,  
775 Suhana Bedi, Nevin Aresh, Joseph Boen, Shriya Reddy, Ben Athiwaratkun, Shuaiwen Leon  
776 Song, and James Zou. Disentangling reasoning and knowledge in medical large language  
777 models. *ArXiv*, abs/2505.11462, 2025. URL [https://api.semanticscholar.org/  
778 CorpusID:278714970](https://api.semanticscholar.org/CorpusID:278714970).
- 779 Andrew Turpin, Jason Wei, Denny Zhou, Quoc V Le, and Ed H Chi. Faithful chain-of-thought  
780 reasoning. *arXiv preprint arXiv:2305.15020*, 2023.  
781
- 782 David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle  
783 Lo, Tom Hope, Luca Soldaini, Shannon Zejiang Shen, Doug Downey, Hannaneh Hajishirzi, and  
784 Arman Cohan. Sciriff: A resource to enhance language model instruction-following over scien-  
785 tific literature, 2024a. URL <https://arxiv.org/abs/2406.07835>.
- 786 David Wadden, Kejian Shi, Jacob Morrison, Aakanksha Naik, Shruti Singh, Nitzan Barzilay, Kyle  
787 Lo, Tom Hope, Luca Soldaini, Shannon Zejiang Shen, et al. Sciriff: A resource to enhance  
788 language model instruction-following over scientific literature. *arXiv preprint arXiv:2406.07835*,  
789 2024b.
- 790
- 791 Changyue Wang, Weihang Su, Qingyao Ai, Yujia Zhou, and Yiqun Liu. Decoupling reasoning and  
792 knowledge injection for in-context knowledge editing, 2025. URL [https://arxiv.org/  
793 abs/2506.00536](https://arxiv.org/abs/2506.00536).
- 794 Pengfei Wang et al. Scienceqa: A large-scale open dataset for question answering in science educa-  
795 tion. *arXiv preprint arXiv:2210.08127*, 2023a.
- 796
- 797 Weijie Wang, Xiang Chen, et al. Evaluating the faithfulness of chain-of-thought reasoning in large  
798 language models. *arXiv preprint arXiv:2401.02392*, 2024a.
- 799 Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R  
800 Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level sci-  
801 entific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*,  
802 2023b.
- 803
- 804 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming  
805 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi  
806 Fan, Xiang Yue, and Wenhao Chen. Mmlu-pro: A more robust and challenging multi-task language  
807 understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024b. URL [https://arxiv.  
808 org/abs/2406.01574](https://arxiv.org/abs/2406.01574).
- 809 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.

- 810 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc  
811 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models,  
812 2023. URL <https://arxiv.org/abs/2201.11903>.
- 813  
814 Juncheng Wu, Sheng Liu, Haoqin Tu, Hang Yu, Xiaoke Huang, James Zou, Cihang Xie, and Yuyin  
815 Zhou. Knowledge or reasoning? a close look at how llms think across domains, 2025. URL  
816 <https://arxiv.org/abs/2506.02126>.
- 817 Xin Xu, Qiyun Xu, Tong Xiao, Tianhao Chen, Yuchen Yan, Jiaxin Zhang, Shizhe Diao, Can Yang,  
818 and Yang Wang. Ugphysics: A comprehensive benchmark for undergraduate physics reasoning  
819 with large language models, 2025. URL <https://arxiv.org/abs/2502.00334>.
- 820 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
821 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint*  
822 *arXiv:2412.15115*, 2024.
- 823  
824 Xiao-Wen Yang, Xuan-Yi Zhu, Wen-Da Wei, Ding-Chu Zhang, Jie-Jing Shao, Zhi Zhou, Lan-Zhe  
825 Guo, and Yu-Feng Li. Step back to leap forward: Self-backtracking for boosting reasoning of  
826 language models, 2025.
- 827 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik  
828 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.
- 829  
830 Ge Zhang et al. Sciglm: Pre-training generalist language models for science with scientific papers.  
831 *arXiv preprint arXiv:2402.00730*, 2024.
- 832  
833 Wayne Xin Zhao et al. A survey of llms for scientific research. *arXiv preprint arXiv:2307.07927*,  
834 2023.
- 835 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
836 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
837 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- 838  
839

## 840 A LIMITATIONS

841  
842 Our KRUX framework and KI extraction methods depend on strong models like DeepSeek-R1 for  
843 generating reasoning traces. While we used an open-weight model, which provides more trans-  
844 parency and interpretability, the KI extraction pipeline may introduce unobservable biases (though  
845 risk is minimal due to our focus on scientific domains), unwanted leakage of information about the  
846 answer, or inconsistencies in the faithfulness of the KIs to the task. To mitigate this, we conducted  
847 manual analysis of the KIs, confirming their relevance and no direct answer leakage, but extracted  
848 KIs could occasionally be irrelevant or incomplete, especially if deployed at scale. Furthermore,  
849 some of our analyses are confounded by factors such as context sensitivity (addressed via permuta-  
850 tions) and the impact of constraining the search space when providing KIs, which we interpret as an  
851 upper bound but may overestimate pure recall benefits. We have taken measures to mitigate these  
852 and discussed the caveats in our discussion of results with more details.

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

859 The benchmarks we examine emphasize STEM fields, which may underrepresent interdisciplinary  
860 or emerging scientific research areas. We acknowledge potential data contamination issues that  
861 may impact our analysis; however, the nature of our study is analytical, and we perform controlled  
862 experiments. In our benchmarks, we also mitigate these concerns by focusing on recent 2024–2025  
863 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	o3			o3-mini			o4-mini			Gemini-2.5-Pro			Claude-Sonnet-4			GPT-5		
	Low	High	$\Delta$	Low	High	$\Delta$	Low	High	$\Delta$	Low	High	$\Delta$	Low	High	$\Delta$	Low	High	$\Delta$
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
SciRIF*	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
<b>Average</b>	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
<b>0.01\$ / Instance</b>	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.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

## C.1 EVALUATION SUITE CURATION

**Subtask Selection Protocol** Our selection process involves joint decisions from three authors, where the two authors determine targeted tasks together and resolve disagreements by discussion, and another author validates them independently, following the exact selection policy. We outline the main criteria of our selection process below:

1. Most candidate benchmarks we considered have clear documentation in their associated manuscripts that describes the characteristics of each subtask, which are indicative of their difficulty and reasoning intensity. For example, SciKnowEval specifies five difficulty levels for the subtasks: L1: Knowledge Memory; L2: Knowledge Comprehension; L3: Knowledge Comprehension; L4: Knowledge Discernment; L5: Knowledge Application. We relied on these descriptions to determine which subtasks to include. For example, for SciKnowEval we only kept subtasks that are labeled with L3 Knowledge Comprehension and above.
2. We then follow the two-step exclusion policy: 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 (i.e., tasks that are not genuinely reasoning-intensive).  
For example, SciEval evaluates from 4 dimensions: basic knowledge, knowledge application, scientific calculation, and research ability. The basic knowledge subset is excluded in the Step 1 paper inspection phase as it demands limited reasoning capabilities by design. After the exclusion inspection, we only keep the knowledge application and scientific calculation subsets, as the research ability subset contains questions that are research-relevant but not reasoning-intensive.
3. Meanwhile, we also exclude tasks with subjects from niche areas that lie outside mainstream STEM (e.g., weapon science, textile engineering from SuperGPQA).

To validate our filtering results, we have another author independently conduct the filtering process from Step 2 again, following the exact same selection policy and compared with the previous results. Among 111 candidate subtasks/domains, SciReas settled on 75 subtasks (listed in Table 13-14), and the two-fold annotation validation reached an agreement accuracy of 90.1%. The disagreement was mainly due to the scope decision in Step 3 for SuperGPQA where tasks from more domains could be within our scope. For example, fields of “Stomatology” and “Public Health and Preventive Medicine” could be justified as mainstream and contain instances that are scientific reasoning-intensive. For SciReas in our submission, we apply a strict filter and exclude fields like these. We list the selection of each benchmark as follows. See Table 13-14 for domain distribution.

**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 13-14) across 3 tasks. License: ODC-BY.

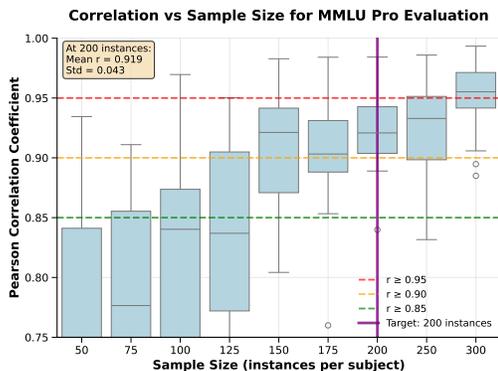


Figure 6: Correlation between sampled and full dataset performance as a function of sample size. 200 instances per subject (purple) yields  $r = 0.919 \pm 0.043$ . Error bars: SD over 30 samples.

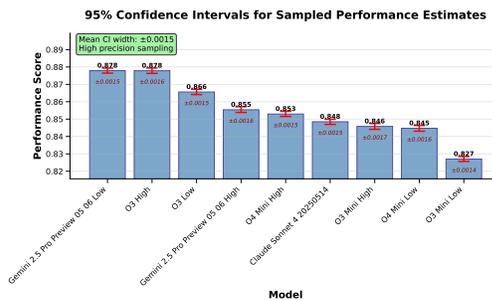


Figure 7: 95% confidence intervals for 200-instance sampling across nine frontier models. Mean CI half-width  $\approx 0.0015$ ; numbers above/below bars show mean and half-width.

**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.

**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

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 (`high_c`, `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.

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.

1080	<b>SYSTEM MESSAGE</b>
1081	
1082	You are an expert judge comparing reasoning intensity between two questions. Analyze both questions thoroughly and determine which one demands more complex reasoning.
1083	Reply in this exact format:
1084	###EXPLANATION: <detailed analysis of both questions and the comparison>
1085	###RESULTS: A / B / UNCLEAR
1086	
1087	
1088	<b>USER MESSAGE</b>
1089	
1090	You will be shown two questions (A and B) from the same academic domain.
1091	A question is <i>*reasoning intensive*</i> if it requires:
1092	• Complex multi-step logical reasoning
1093	• Advanced mathematical computation or derivation
1094	• Integration of multiple concepts or principles
1095	• Abstract thinking or sophisticated problem-solving strategies
1096	• Deep domain knowledge application
1097	<i>*QUESTION A*</i>
1098	Context: {{context_a}}
1099	Question: {{question_a}}
1100	<i>*QUESTION B*</i>
1101	Context: {{context_b}}
1102	Question: {{question_b}}
1103	Analyze both questions carefully and explain your reasoning. Then reply using the exact format specified above.

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

1107	<b>SYSTEM MESSAGE</b>
1108	
1109	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.
1110	Reply in this exact format:
1111	###EXPLANATION: <detailed analysis of why the low-reasoning model failed>
1112	###RESULTS: REASONING/KNOWLEDGE/BOTH/UNCLEAR
1113	
1114	
1115	
1116	
1117	<b>USER MESSAGE</b>
1118	
1119	You will be shown a question from an academic dataset, along with
1120	(1) a <i>*CORRECT*</i> answer from a high-reasoning model and
1121	(2) an <i>*INCORRECT*</i> answer from a low-reasoning model.
1122	Your task is to analyze <i>*why*</i> the low-reasoning model failed.
1123	Consider whether the failure is primarily due to:
1124	• <i>*REASONING*</i> : Insufficient logical thinking, problem-solving ability, or step-by-step analysis
1125	• <i>*KNOWLEDGE*</i> : Lack of domain knowledge (missing facts, formulas, concepts, procedures)
1126	• <i>*BOTH*</i> : Significant deficiencies in both reasoning and knowledge
1127	• <i>*UNCLEAR*</i> : Cannot determine the primary cause of failure
1128	<b>QUESTION</b>
1129	Context: {{context}}
1130	Question: {{question}}
1131	<b>CORRECT ANSWER (from {{high_model}}):</b>
1132	{{high_full_response}}
1133	<b>INCORRECT ANSWER (from {{low_model}}):</b>
	{{low_full_response}}
	Analyze why the low-reasoning model failed. Was it primarily due to insufficient reasoning ability or lack of knowledge?

Figure 9: Prompt used to classify failure cause (reasoning vs. knowledge) for low-reasoning models.

**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,<sup>7</sup> 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,<sup>8</sup> and top-p=0.95, borrowed from recommended setting by Anthropic,<sup>9</sup> 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)
<i>OpenAI models</i>		
GPT-4.1-2025-04-14	2.00	8.00
o3-mini-2025-01-31	1.10	4.40
o3-2025-04-16	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
<i>DeepSeek models</i>		
DeepSeek-V3-0324	0.14	0.28
DeepSeek-R1-0120	0.55	2.19
DeepSeek-R1-0528	0.55	2.19
<i>Alibaba Qwen models (Together AI)</i>		
Qwen3-32B	0.40	1.20
Qwen3-235B-2507	0.65	3.00
<i>Google models</i>		
Gemini-2.5-Pro-Preview-05-06	1.25	10.00
<i>Meta models (Together AI)</i>		
Llama-4-Maverick-17B-128E-Instruct-FP8	0.27	0.85
<i>Anthropic models</i>		
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  $f_\theta$  to approximate a distribution over output

<sup>7</sup>[https://docs.litellm.ai/docs/providers/anthropic#usage—thinking—reasoning\\_content](https://docs.litellm.ai/docs/providers/anthropic#usage—thinking—reasoning_content)

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

<sup>9</sup><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:  $\mathcal{L}(\theta) = -\sum_{(x,y) \in \mathcal{D}} \sum_{t=1}^{|y|} \log p_{\theta}(y_t | y_{<t}, x)$ , 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.<sup>10</sup> 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., “cm<sup>2</sup>”) or chemical formulas (e.g., “H<sub>2</sub>O” 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. Qwen-Math slightly outperforms Qwen-STEM on SciREAS and the gap is amplified on SciREAS-PRO.

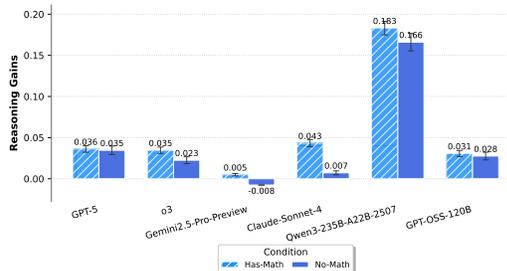


Figure 10: Performance gains on Has-Math instances vs. No-Math instances across different frontier models when reasoning effort increases. CI95% shown as the error bar. CoT helps more with Has-Math instances.

<sup>10</sup>Longer input lengths would slow down our training in quadratic order based on 8 80GB A100/H100 GPUs.

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

Units recognised by the heuristic		
• % • °C, °F, K, • g, kg, mg, • μg/ug, lb/lbs, • oz	• J, kJ, MJ; W, • kW, MW, GW • V, kV; A, mA, • μA/uA	• mol; M, mM, • μM/uM, nM, • pM • dB; rpm; • rad/s
• m, cm, mm, • km • L/l, mL/ml, • μL/μl/ul	• Hz, kHz, • MHz, GHz • cm <sup>2</sup> , m <sup>2</sup> , • mm <sup>2</sup> , km <sup>2</sup> , • cm <sup>3</sup> , m <sup>3</sup> , • mm <sup>3</sup> , km <sup>3</sup>	• s, ms, μs/us, • ns; min, h • day/days; • yr/yrs
• Pa, kPa, MPa, • atm, bar, mbar		

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

Given limited subject coverage on SYNTHETIC-1-Math dataset, the strong performance of checkpoints 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.

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 non-math instances for SCIREAS-PRO. -Math and -STEM variants contribute to different dimensions of performance, while -BOTH captures improvements on both.

Models	OpenR1	Llama-Nemotron	General-Reasoner	SYNTHETIC-1-SFT	Qwen-Nemotron	Qwen-BOTH	SciLIT01	Qwen3	Qwen3-thinking
	Q2.5-Math	L3.1-Inst.	Q2.5-Base	Q2.5-Inst.			Q3-Base		
Training Methods Trained by Us	SFT No	SFT&RL No	RL No	SFT No	SFT Yes	SFT Yes	SFT Yes	- No	- No
GPQA	44.42	37.95	35.94	38.84	44.20	40.63	50.89	55.80	55.80
SuperGPQA*	31.90	29.39	14.26	22.39	19.47	20.33	30.11	23.32	38.27
MMLU-Pro*	60.86	65.64	62.14	56.21	63.57	65.00	76.57	73.36	81.71
LabBench*	27.14	27.78	35.58	28.61	35.76	33.00	35.07	36.99	38.19
OlympiadBench	53.03	37.62	19.82	40.75	29.33	34.55	43.78	28.51	21.30
SciBench	61.85	57.66	19.08	51.59	48.27	47.11	61.27	54.05	68.21
SciEval*	43.64	68.67	70.34	46.41	38.53	72.36	80.60	81.51	84.02
SciKnowEval*	28.45	30.69	34.19	19.13	31.85	32.00	39.46	37.99	41.81
SciRIFF*	29.17	34.01	37.75	28.57	39.24	41.81	44.01	47.23	47.26
UGPhysics*	50.30	45.92	20.86	43.96	46.52	40.03	52.28	30.98	59.81
<b>Average</b>	43.08	43.53	34.99	37.64	39.67	42.68	<u>51.41</u>	46.97	<b>53.64</b>
<b>SciREAS-PRO</b>	<u>26.43</u>	23.75	13.73	19.44	19.68	21.11	24.84	19.05	<b>29.92</b>

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.

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

```

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 gener-
alized 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:
1. 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:
1. Morpholino oligonucleotides can bind to the 5' splice site of pre-mRNA to block spliceosome recog-
nition.
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.
4. Antisense molecules function by complementary base-pairing to specific RNA sequences to modu-
late splicing.
5. Lariat structures are formed during canonical splicing via branch site interactions but are absent
when spliceosome activity is inhibited.
6. R-loops are DNA-RNA hybrid structures unrelated to spliceosome-mediated exon skipping mecha-
nisms.
7. 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.

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

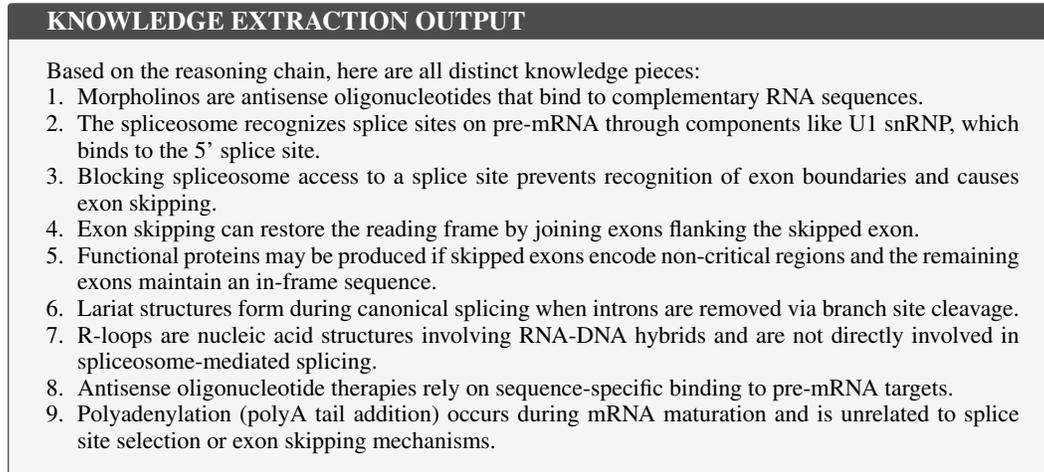


Figure 15: Example KIs extracted from DeepSeek-R1.

## F.2 ALTERNATIVE EXTRACTOR

### F.2.1 QUALITY INSPECTION

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, where one could contain more details than another. For example, given the same CoT traces on a math problem, DeepSeek-R1 and Qwen3 generate KIs as shown in Figure 16. Notably, the two sets of KIs **cover nearly identical problem-specific knowledge**, although DeepSeek phrases its conversion rule abstractly (in terms of quantities, rather than Qwen’s problem-specific cups and gallons) and Qwen3 includes a problem-solving heuristic to avoid common errors.

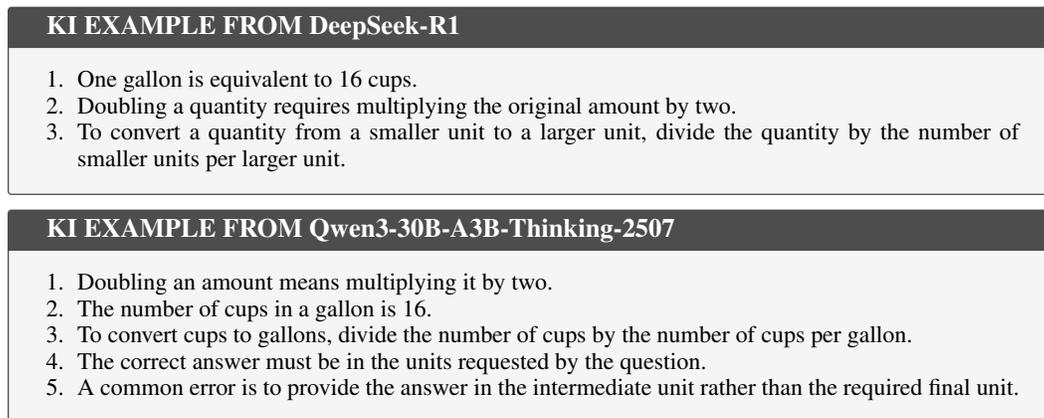


Figure 16: An example of KIs generated by DeepSeek-R1 and Qwen3-30B-A3B-Thinking-2507. The DeepSeek-R1 tends to cover knowledge recalled in more detail.

### F.2.2 USING ALTERNATIVE EXTRACTOR ON KRUX

Running KRUX with KIs extracted by Qwen3-30B-A3B-Thinking-2507 (Qwen3 for short), we present Table 11-12 to show that KIs extracted from DeepSeek-R1 show significant gains compared to the baselines, as our main results in Table 3. With Qwen3 as the extractor, setups with DeepSeek-R1’s KIs in-context (w/ R1 KIs) show significant improvements over the base setups with KIs extracted from the model itself (w/ self KIs) on GPQA in Table 11. A similar trend is shown on MMLU-Pro\* in Table 12 as well. Notably, performance improvements persist in both base models and reasoning-enhanced models. The result shows the same trend as results derived from using R1 as the extractor, demonstrating that KRUX generalizes to different extractor models.

Table 11: Accuracy of Qwen and Llama variants on GPQA with external knowledge ingredients (KIs) extracted by Qwen3-30B-A3B-Thinking-2507 or DeepSeek-R1 as the extractor. w/ R1 KIs setups outperform w/ self KIs setups across different models.

Extractor Setup	Qwen3			DeepSeek-R1		
	w/ self KIs	w/ R1 KIs	$\Delta$	w/ self KIs	w/ R1 KIs	$\Delta$
Qwen	35.0	43.5	+8.5	34.2	47.2	+13.0
Qwen-STEM	42.6	51.1	+8.5	41.6	52.5	+10.9
Qwen-MATH	38.8	50.2	+11.4	39.5	53.5	+14.1
Qwen-BOTH	43.1	52.9	+9.8	40.8	54.5	+13.7
Llama	29.0	39.5	+10.5	29.1	43.6	+14.5
Llama-STEM	38.6	51.8	+13.2	39.0	53.2	+14.2
Llama-MATH	33.9	50.0	+16.1	36.2	53.8	+17.6
Llama-BOTH	39.1	52.2	+13.2	39.4	54.7	+15.3

Table 12: Accuracy of Qwen and Llama variants on MMLU-Pro\* with external knowledge ingredients (KIs) extracted by Qwen3-30B-A3B-Thinking-2507 or DeepSeek-R1 as the extractor. w/ R1 KIs setups outperform w/ self KIs setups across different models.

Extractor Setup	Qwen3			DeepSeek-R1		
	w/ self KIs	w/ R1 KIs	$\Delta$	w/ self KIs	w/ R1 KIs	$\Delta$
Qwen	61.7	66.4	+4.6	59.0	68.9	+9.8
Qwen-STEM	63.8	69.9	+6.1	64.7	69.7	+5.0
Qwen-MATH	65.2	73.6	+8.4	66.9	74.0	+7.1
Qwen-BOTH	65.4	71.4	+6.0	65.7	71.6	+5.9
Llama	49.0	56.4	+7.4	47.7	60.5	+12.8
Llama-STEM	58.8	65.7	+6.9	59.1	68.2	+9.1
Llama-MATH	60.7	66.4	+5.7	59.7	69.0	+9.4
Llama-BOTH	62.6	70.2	+7.6	63.8	72.7	+8.9

### 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.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

```

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 proposed 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"]
}
# Knowledge Statement: {src_text}
Please provide your response in the specified format without any additional text.

```

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

```

EXAMPLE 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 GPT-o3 and GPT-5 from OpenAI for grammar and typo corrections.

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

Domain	Task Source	Subtask/Subdomain	Instances	Total	Metrics	
Physics	GPQA	Physics	187	5087	Acc	
	MMLU-Pro	physics	200		Acc	
	SciBench	fund	81		Acc	
		thermo	83		Acc	
	OlympiadBench-COMP	class	63		Acc	
		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
			GeometricalOptics		54	Acc
			Relativity		200	Acc
			ClassicalElectromagnetism		200	Acc
			ClassicalMechanics		200	Acc
			WaveOptics		200	Acc
			QuantumMechanics		200	Acc
			TheoreticalMechanics		200	Acc
			AtomicPhysics		200	Acc
	SemiconductorPhysics				148	Acc
			154	Acc		
StatisticalMechanics			200	Acc		
			1482	Acc		
Chemistry	GPQA	Chemistry	183	2158	Acc	
	MMLU-Pro	chemistry	200		Acc	
	SciBench	quan	41		Acc	
		chemc	47		Acc	
	atkins	matter	121		Acc	
			57		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		F1, LM	
	SuperGPQA	Computer Science and Technology	108		Acc	
Math	MMLU-Pro	math	200	2533	Acc	
		SciBench	calc		52	Acc
	stat	diff	92		Acc	
			55		Acc	
	OlympiadBench-COMP	maths_en	674		Acc	
	SuperGPQA	Mathematics	1460		Acc	

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

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

Domain	Task Source	Subtask	Instances	Total	Metrics
Biology	GPQA	Biology	78	1911	Acc
	MMLU-Pro	biology	200		Acc
	LabBench	CloningScenarios	33		Acc
		ProtocolQA	108		Acc
		SeqQA	600		Acc
	SciKnowEval.L5	biological_procedure_generation	200		LM
		biological_reagent_generation	200		LM
	SciEval	biology_knowledge_application	200		Acc
		biology_scientific_calculation	200		Acc
	SuperGPQA	Biology	92		Acc
Medicine	MMLU-Pro	health	200	634	Acc
	SciRIFB	SciFact	184		F1, LM
		Evidence Inference	250		F1
Material Sci	SciKnowEval.L5	crystal_structure_and_composition	196	624	LM
		specified_band_gap_material_generation	200		LM
		property_and_usage_analysis	118		LM
	SuperGPQA	Materials Science and Engineering	110		Acc
Engineering	MMLU-Pro	engineering	200	2205	Acc
	SuperGPQA	Control Science and Engineering	77		Acc
		Information and Communication Engineering	156		Acc
		Electrical Engineering	234		Acc
		Chemical Engineering and Technology	226		Acc
		Power Engineering and Engineering	345		Acc
		Thermophysics			
		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 14: Domain-wise breakdown of SCIREAS tasks and instance counts (Part 2: Biology to Engineering).

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

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

Table 15: Domain-wise breakdown of SciREAS-PRO tasks and instance counts (Part 1: Physics to Math).

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

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

Table 16: Domain-wise breakdown of SCIREAS-PRO tasks and instance counts (Part 2: Biology to Engineering).