

Demystifying Scientific Problem-Solving in LLMs by Probing Knowledge and Reasoning

Alan Li* Yixin Liu* Arpan Sarkar Doug Downey Arman Cohan

Yale



HARVARD
UNIVERSITY



Northwestern
University



Introduction

- Scientific tasks demand **domain knowledge** *and* **multi-step reasoning**
- Existing benchmarks probe only a **narrow slice** of skill space:
 - Subject coverage (e.g. GPQA - {phys, chem, bio}, LabBench - {bio})
 - MCQ only - {GPQA, MMLU-Pro, LabBench, ...}

How do we evaluate scientific reasoning comprehensively?

How do reasoning and knowledge interact in science problem-solving?

SciREAS: A new benchmark for scientific reasoning

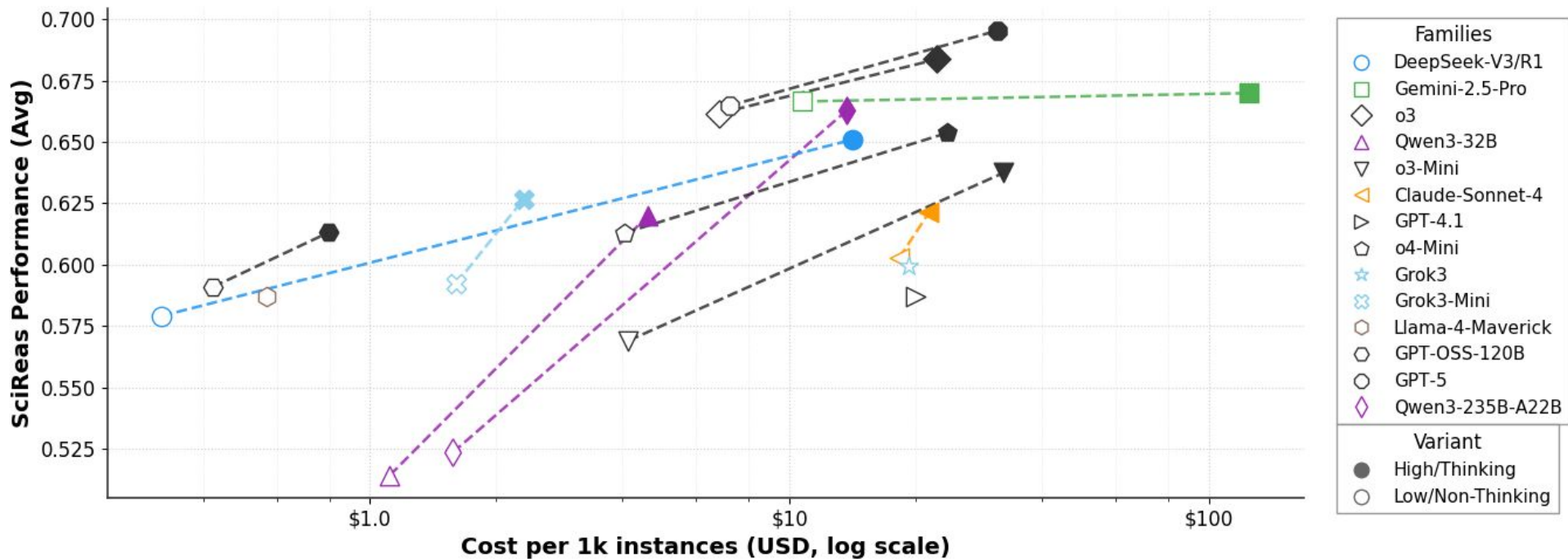
Existing scientific benchmarks: Domain-siloed, fact-heavy, overlook reasoning bottlenecks.

We manually filter existing scientific benchmarks and focus on **reasoning**

(10 datasets)

Sampling to cut the evaluation costs in half, without reducing confidence

□ **SciREAS Benchmark!**



Why holistic evaluation is important?

Performance on subsets do not necessarily correlate

Individual models may have been tuned towards different benchmarks



How does the interplay between **knowledge** and **reasoning** shape scientific problem-solving?

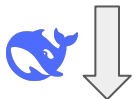
KRUX: Knowledge and Reasoning Utilization Exam

1. Collect reasoning traces using a strong reasoning model
2. Extract **knowledge ingredients K** from reasoning traces (**KI**)
3. Augment testing instances with KI in the prompt

Question Q: While operating on variable frequency supplies, the AC motor requires variable voltage in order to _____. (A) extend the motor's lifespan. (B) increase the motor's efficiency. (C) avoid effect of saturation (D) ...



Reasoning R + Answer A: <think> Okay, so I need to figure out why an AC motor requires variable voltage when operating on variable frequency supplies ... </think>...Therefore, the answer is (C).

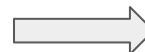


Knowledge Pieces:

KI₁: The synchronous speed of an AC motor is proportional to the ratio of supply frequency to the number of motor poles.

KI₂: Induction motors require maintenance of a constant voltage-to-frequency ratio for optimal operation.

...



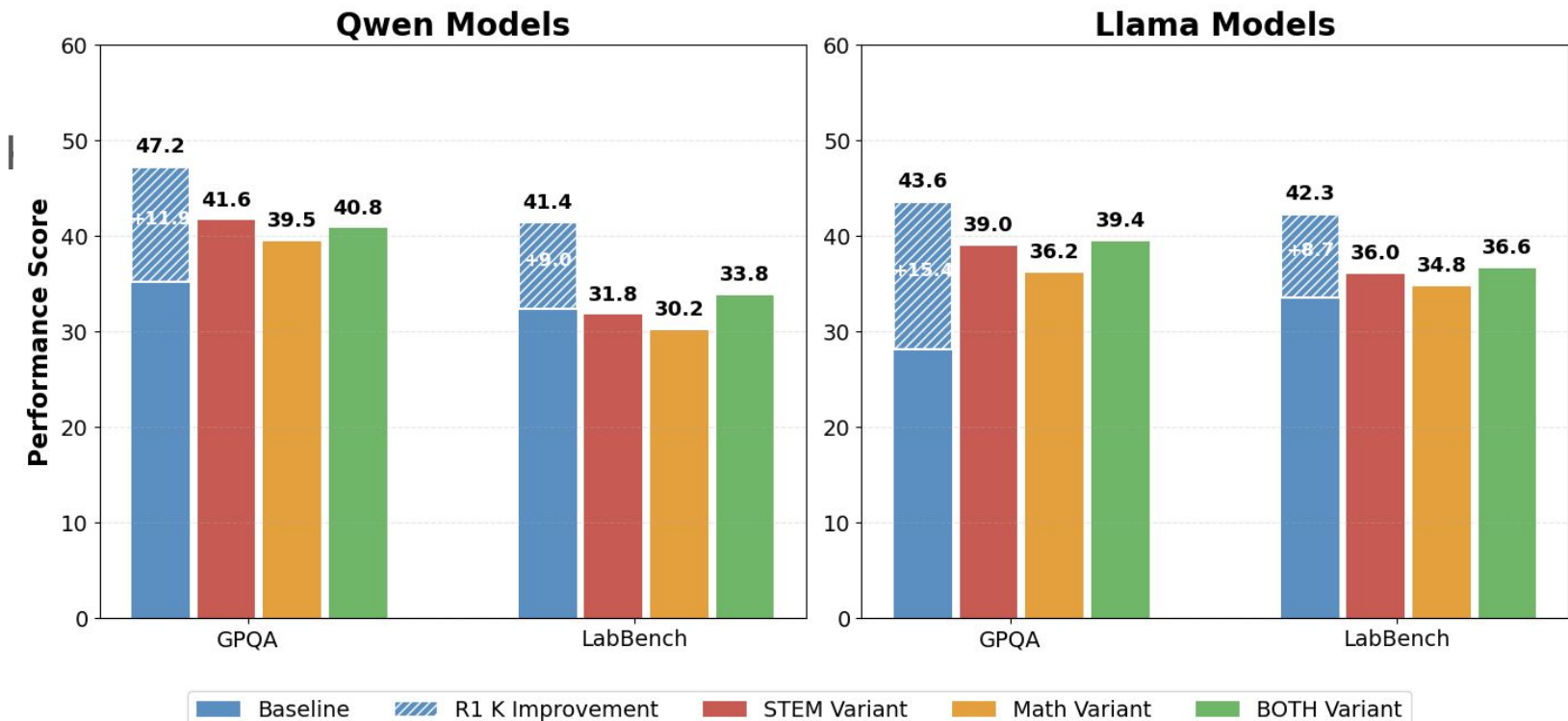
Question w/ KI

Question Q: ...

Here are some knowledge points that could be helpful:

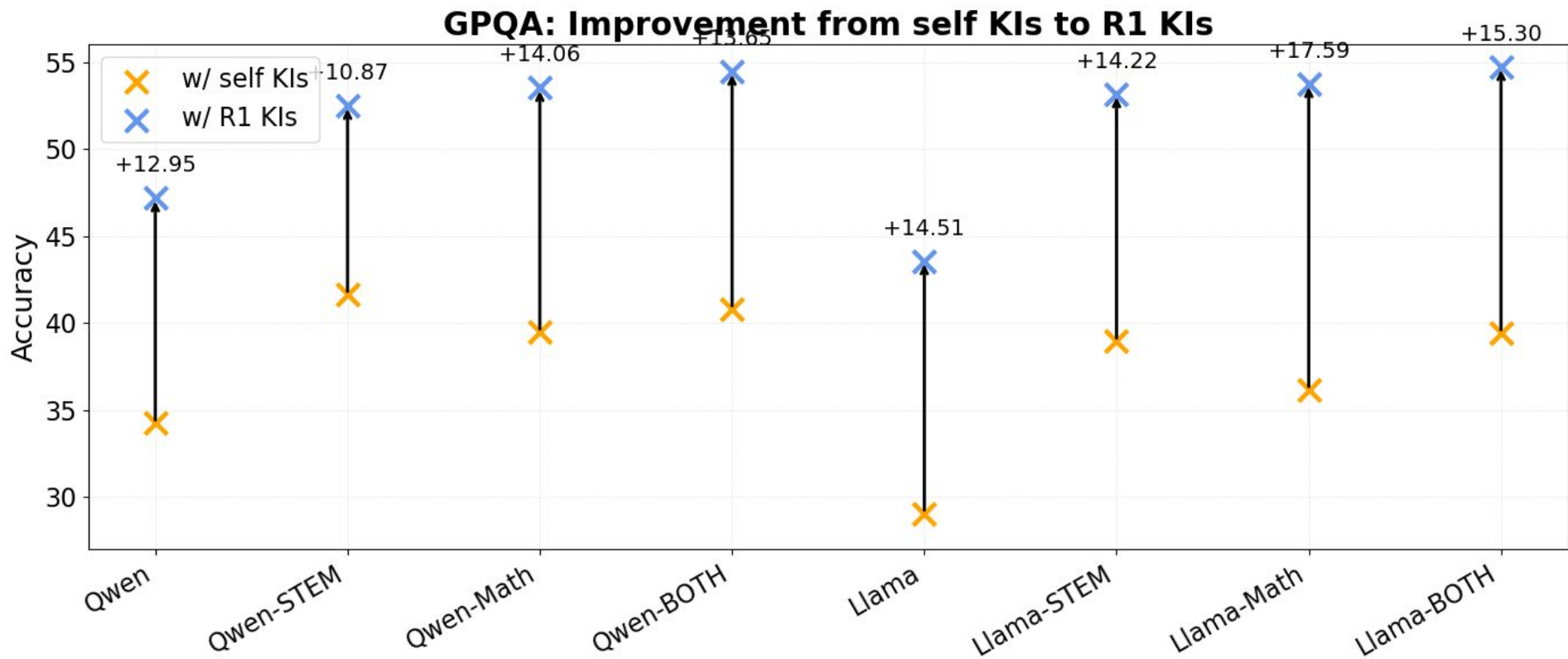
- **KI₁**
- **KI₂**

RQ1: Knowledge Bottleneck



Missing knowledge hinders model performance

RQ2: Additive Gains



KIs extracted from different post-trained models and supplied to the prompt

RQ3: how different KIs affect the same base model

Base Setup		GPQA	MMLU-Pro*
Qwen	w/ Qwen KIs	34.24 \pm 0.93	59.03 \pm 0.34
	w/ Qwen-Math KIs	36.93 \pm 1.75	63.66 \pm 0.45
Llama	w/ Llama KIs	29.06 \pm 1.44	47.73 \pm 0.89
	w/ Llama-Math KIs	29.69 \pm 1.72	53.91 \pm 0.94

Findings: Reasoning-focused training on math tends to help the model in surfacing the knowledge better.

Takeaways

Need for comprehensive benchmarking in scientific problem solving

Models could be bottlenecked by lack of specialized knowledge

Reasoning models consistently prevail with the use of context

Reasoning-focused post-training helps models in surfacing and utilizing knowledge

Thanks!