



Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?

Yang Yue^{1*†}, Zhiqi Chen^{1*}, Rui Lu¹, Andrew Zhao¹, Zhaokai Wang², Yang Yue¹, Shiji Song¹, and Gao Huang^{1‡}

¹ Tsinghua University ² Shanghai Jiao Tong University
* Equal Contribution † Project Lead ‡ Corresponding Author

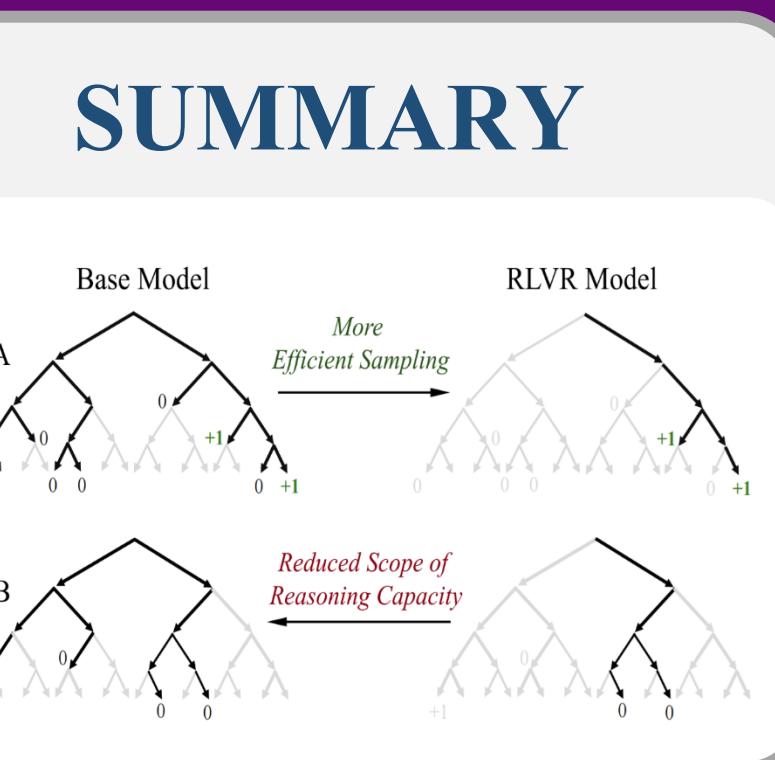
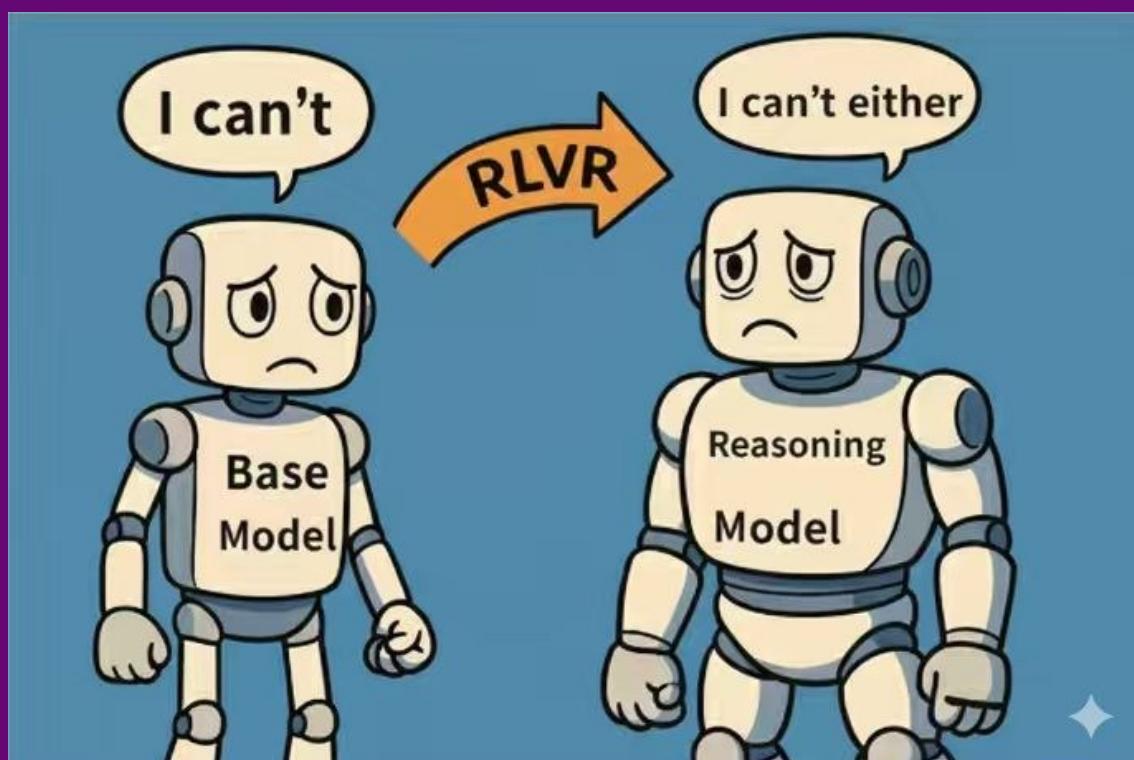


[paper](#) [homepage](#)

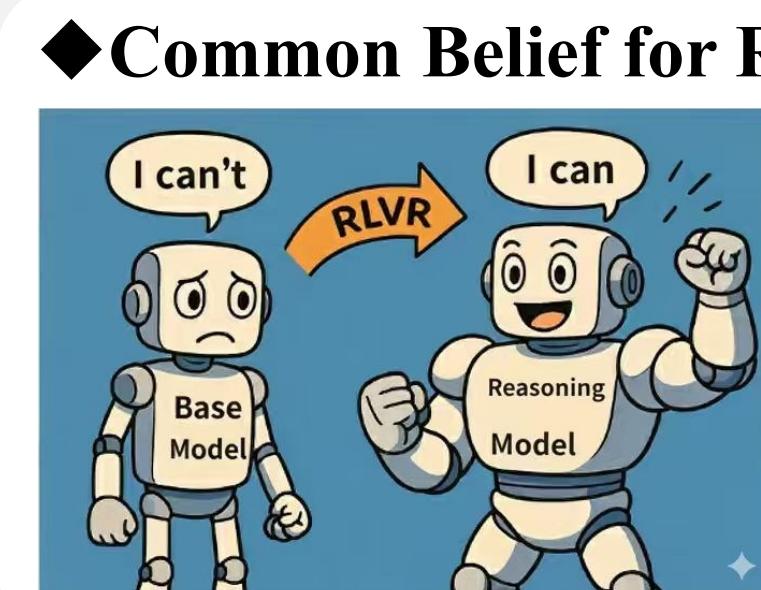
oral



RLVR rarely expands reasoning boundary, mainly improves sampling efficiency!

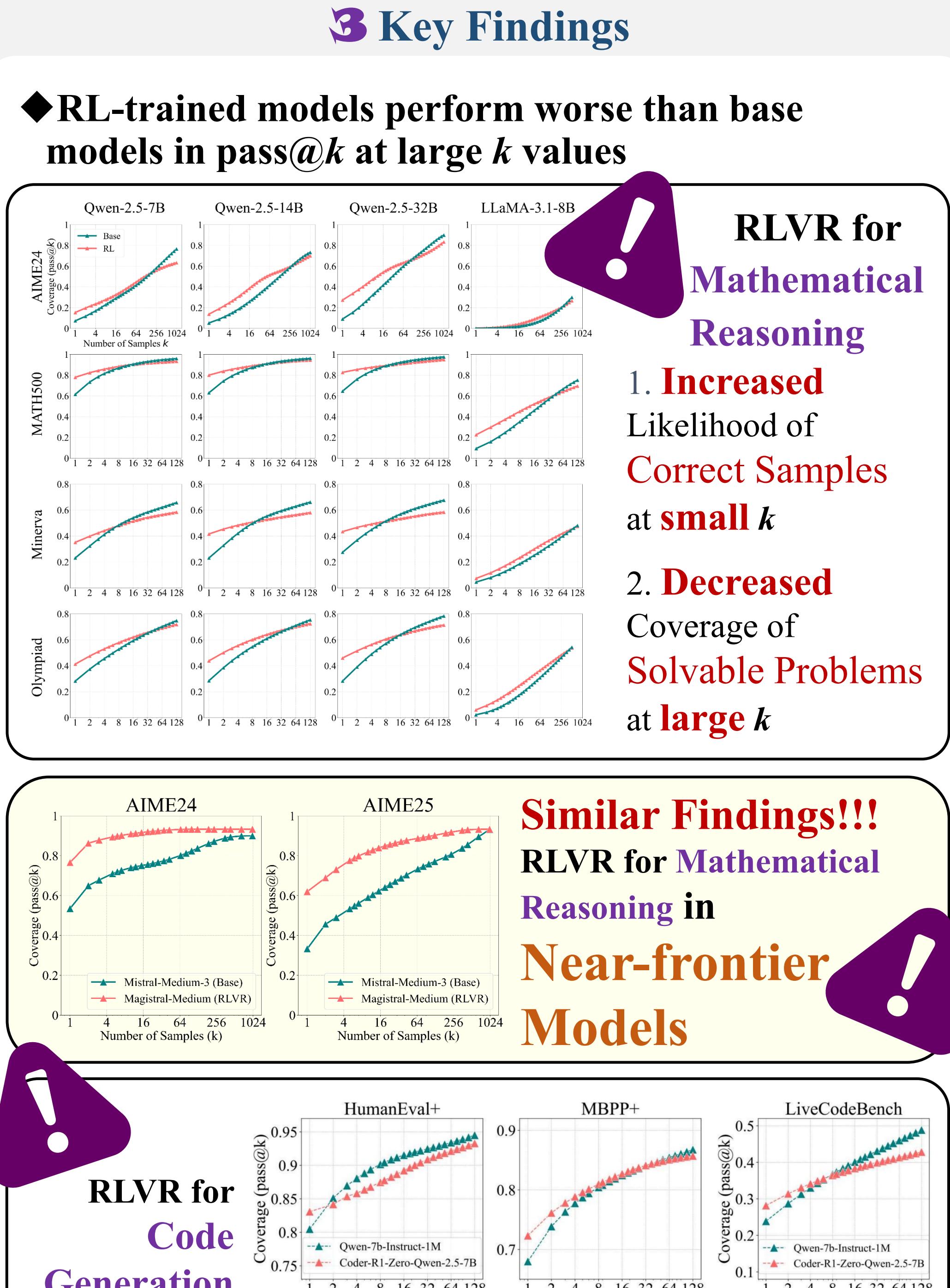
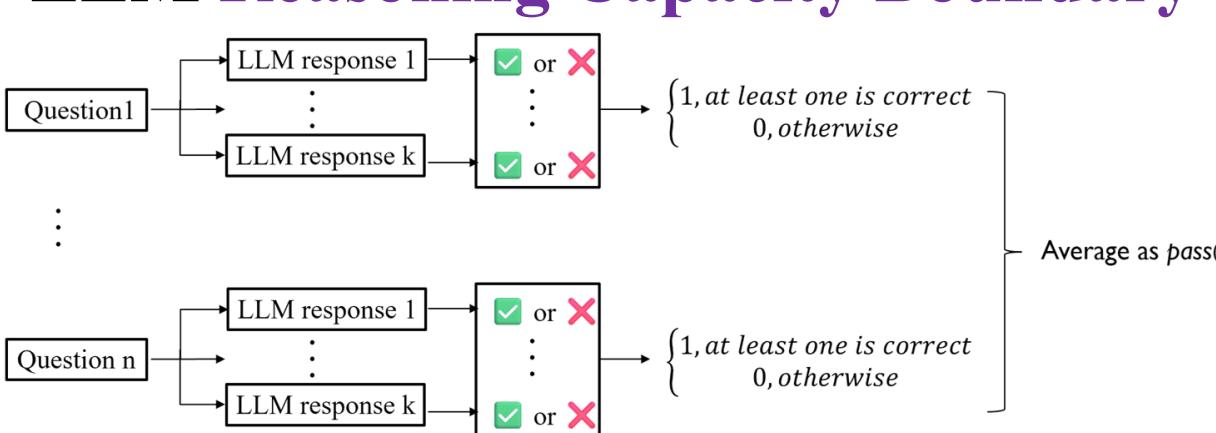


1 Background & Preliminaries



◆ Common Belief for RLVR

◆ [pass@k] A metrics for LLM Reasoning Capacity Boundary



2 Experiment Setup

Task	Start Model	RL Framework	RL Algorithm(s)	Benchmark(s)
Mathematics	LLaMA-3.1-8B Qwen2.5-7B/14B/32B-Base Qwen2.5-Math-TB	SimpleRLZoo Oat-Zero	GRPO	GSM8K, MATH500 Minerva, Olympiad AIME24, AMC23 LiveCodeBench
Code Generation	Qwen2.5-7B-Instruct	Code-R1	GRPO	HumanEval+ Coder-RL-Zero-Qwen-2.5-7B
Visual Reasoning	Qwen2.5-7B-Base Qwen2.5-7B-Instruct	PPO	Reinforce++	MBPP+
Deep Analysis	DeepSeek-R1-Distill-Qwen-7B	VeRL	RLOO, ReMax, DAPO	LiveCodeBench HumanEval+ MathVista MathVision Omni-Math-Rule MATH500

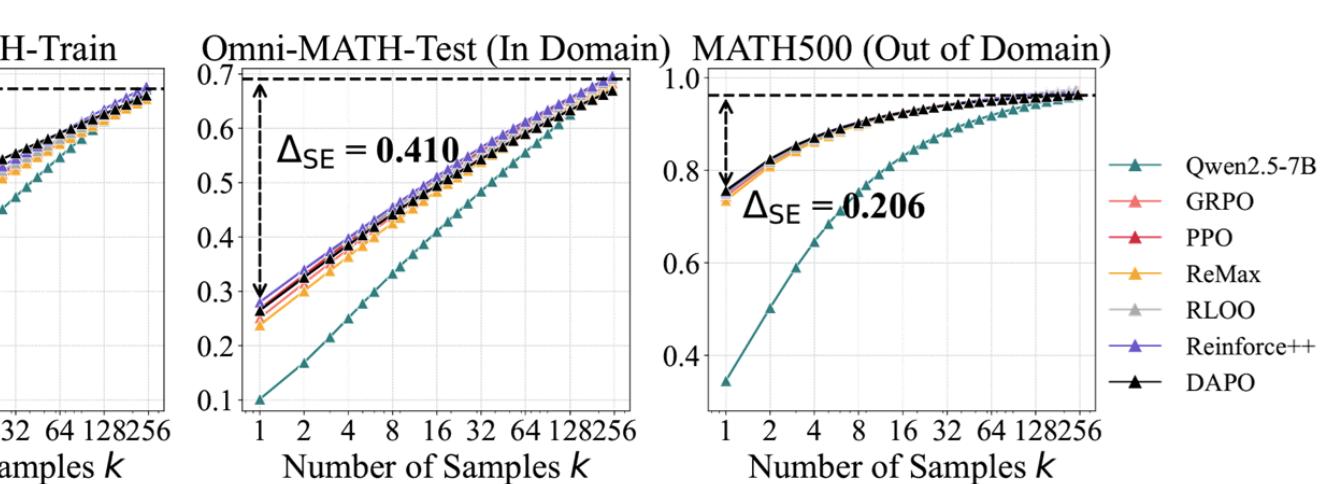
Keep prompt template for base model identical to RLVR template

4 Deep Analysis

◆ RLVR algorithms

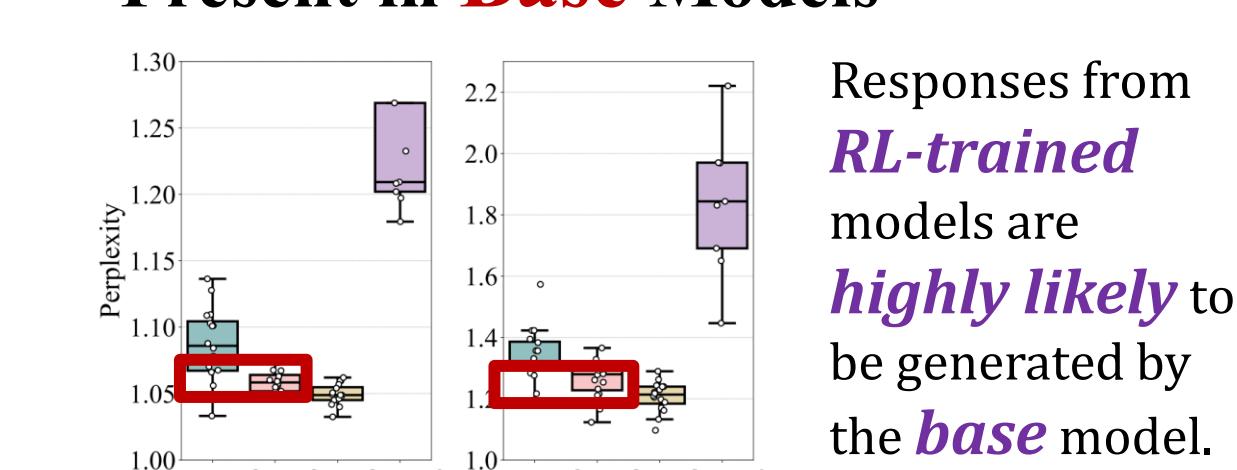
perform similarly and remain far from optimal

Define $\Delta_{SE} = \text{Base pass}@k - \text{RL pass}@1$.



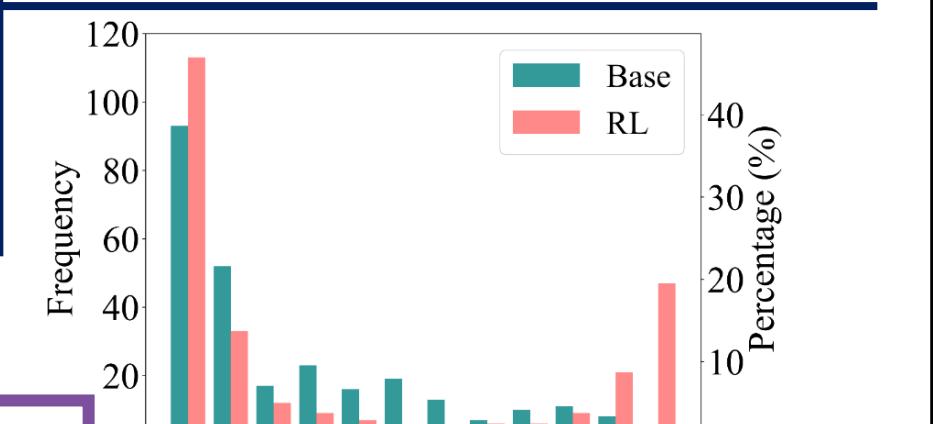
◆ As RL training progresses, pass@1 consistently improves while pass@256 decreases

◆ Reasoning Patterns Maybe Already Present in Base Models



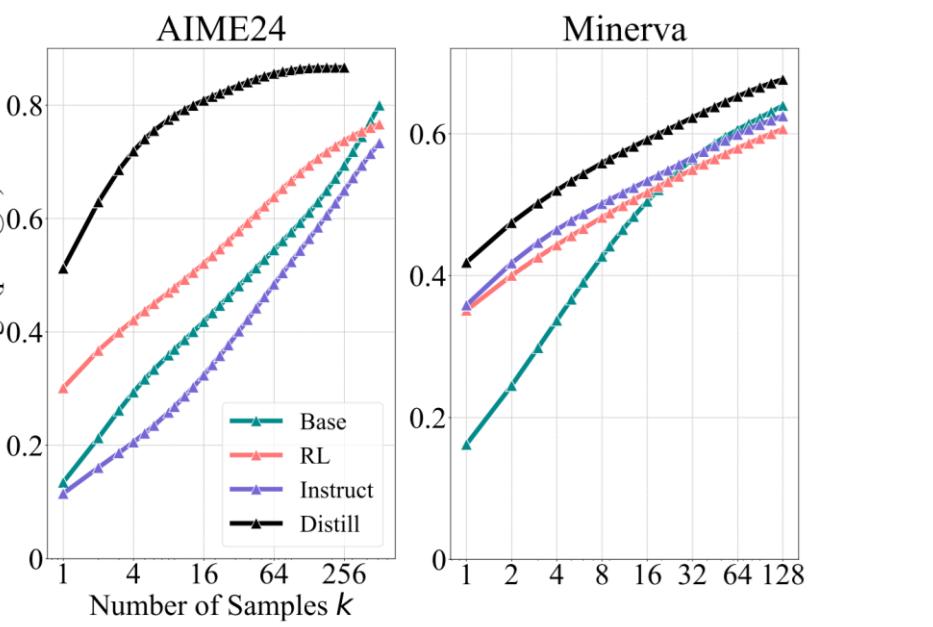
Base	SimpleRLZoo	AIME24	MATH500
✓	✓	63.3%	92.4%
✓	✗	13.3%	3.6%
✗	✓	0.0%	1.0%
✗	✗	23.3%	3.0%

The RLVR Model's Solvable-Problem Coverage is an *Approximate Subset* of the Base Model's Coverage.



The improvement in average scores, driven by improving sampling efficiency on already solvable problems.

◆ RLVR and distillation are fundamentally different: distillation expands reasoning by injecting new knowledge.



Potential Future Directions

1. RL data & compute scale up!
2. Agent with tools & external info
3. Exploration mechanism (AlphaEvolve)
4. Process reward & better credit assignment

