# Reinforcement Learning with Verifiable Rewards

Madhav Donepudi

Tanishk Govil

David Jagga

# Background

RLVR: RL where reward is **automatically checkable**

Why people care: offers a scalable way to push LLMs towards **correct** outcomes without expensive human grading

Core Pipeline: Prompt → LLM output → Verifier

Three papers:

1. **R1:** RLVR yields strong reasoning gains
2. **Boundary Paper:** asks whether gains reflect efficiency vs new capacity
3. **ProRL:** shows conditions where longer RL can expand boundaries

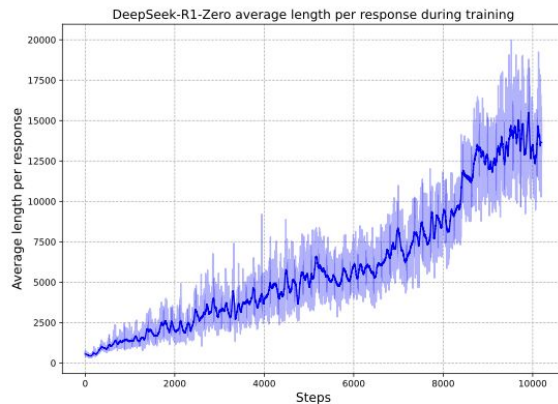# DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning
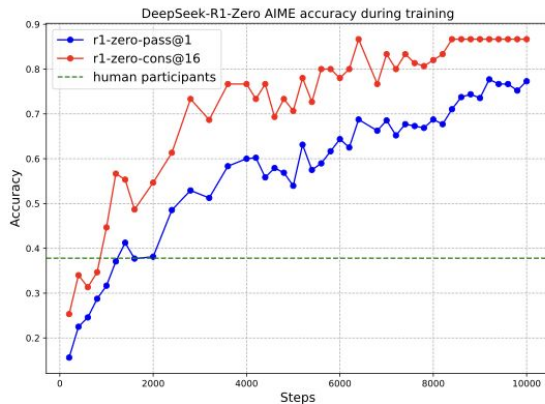
DeepSeek-AI, 2025
Guo et al., arXiv:2501.12948

# The Bottleneck of Supervised Fine-Tuning

- Problem: SFT relies on expensive human data and caps performance at human levels

- Research Question: Can reasoning emerge purely from Reinforcement Learning without human-labeled trajectories?

- Hypothesis: Scaling "test-time compute" via RL allows models to discover superior, non-human reasoning paths
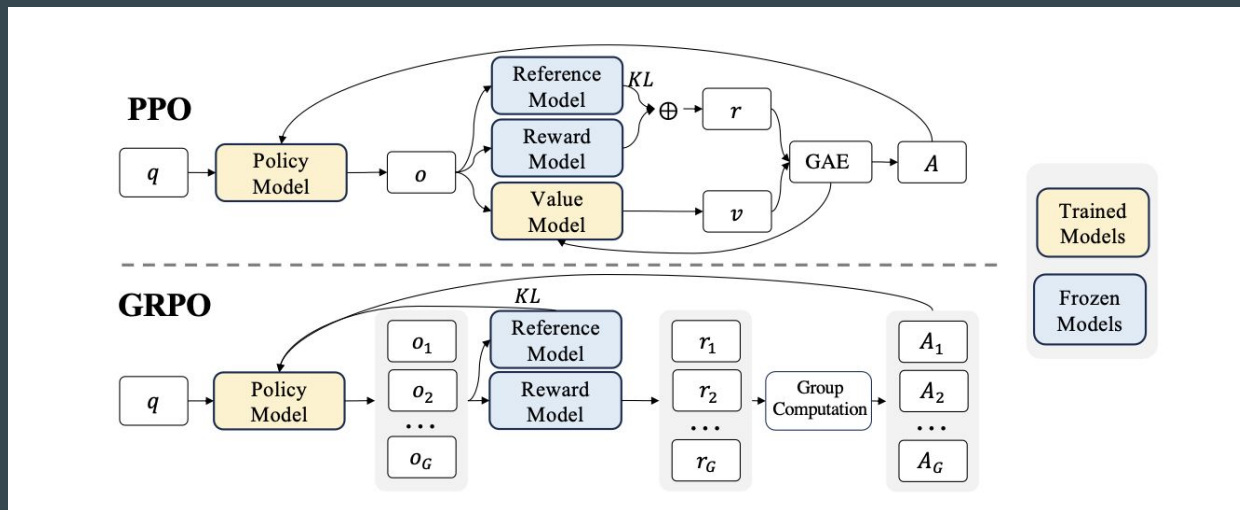
# DeepSeek-R1-Zero: Pure RL Discovery

- Initialization: Trained directly from DeepSeek-V3-Base with Zero SFT
- Reward Structure: Solely rule based (Accuracy + Format)
- Accuracy:
  - Math: Correct answer presented in the correct format (eg. within a box)
  - Code: Compiler verifies if code passes specific test cases
- Format:
  - Structural Check: Encourages thinking process within mandatory <think> and <answer> tags
- Optimization: Pure exploration and no human style guidance

# Group Relative Policy Optimization (GRPO)

- Motivation: Standard PPO is too memory intensive due to the "Value Model"
- Innovation: Advantages are computed relative to a group of outputs
- Impact: Significant VRAM savings which enables large scale training of long-chain reasoning

# Mathematical Framework of GRPO

- Advantage: Standardized reward within a group of G outputs using advantage normalization
  - Correct response: Non-negative Advantage
  - Incorrect response: Non-positive Advantage
- KL Penalty: Unbiased estimator used to maintain policy stability

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^{G} \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G}\sum_{i=1}^{G}\left(\min\left(\frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}A_i, \text{clip}\left(\frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1-\varepsilon, 1+\varepsilon\right)A_i\right) - \beta\mathbb{D}_{KL}\left(\pi_\theta||\pi_{ref}\right)\right), \quad (1)$$
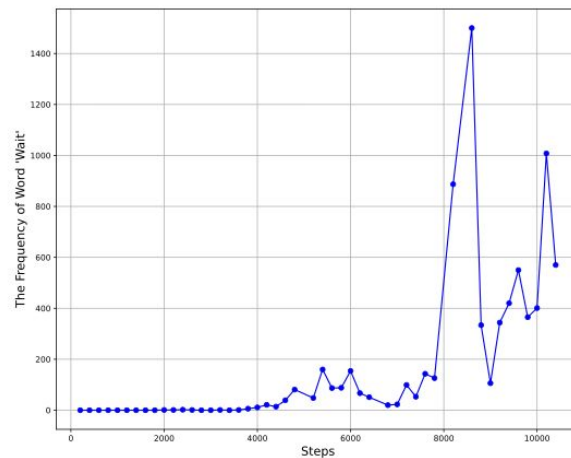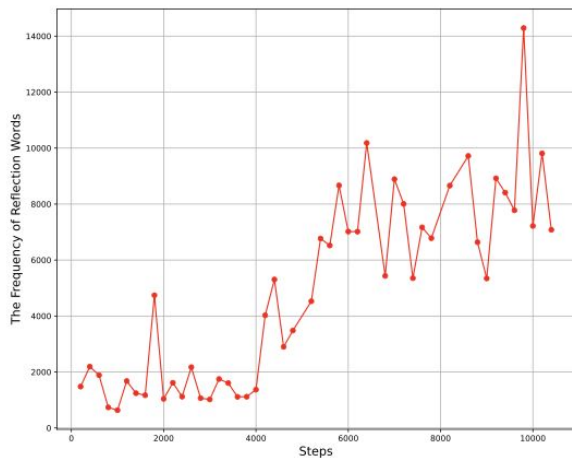
$$\mathbb{D}_{KL}\left(\pi_\theta||\pi_{ref}\right) = \frac{\pi_{ref}(o_i|q)}{\pi_\theta(o_i|q)} - \log\frac{\pi_{ref}(o_i|q)}{\pi_\theta(o_i|q)} - 1, \quad (2)$$

$$A_i = \frac{r_i - \text{mean}(\{r_1, r_2, \cdots, r_G\})}{std(\{r_1, r_2, \cdots, r_G\})}. \quad (3)$$

1. GRPO Objective     2. KL Penalty Term     3. Advantage from verifier rewards

# Emergent Self-Correction and Scaling

- Discovery: Model autonomously learned to backtrack, verify, and re-evaluate
- Linguistic Markers: Sharp increase in "reflective" tokens (e.g., "Wait", "Let me re-check")

# The "Aha Moment"

Table 2 | An interesting "aha moment" of an intermediate version of DeepSeek-R1-Zero. The model learns to rethink using an anthropomorphic tone. This is also an aha moment for us, allowing us to witness the power and beauty of reinforcement learning.

---

**Question:** If $a > 1$, then the sum of the real solutions of $\sqrt{a - \sqrt{a + x}} = x$ is equal to

---

**Response:** `<think>`

To solve the equation $\sqrt{a - \sqrt{a + x}} = x$, let's start by squaring both $\cdots$

$$\left(\sqrt{a - \sqrt{a + x}}\right)^2 = x^2 \implies a - \sqrt{a + x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

$\cdots$

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be $\cdots$

We started with the equation:

$$\sqrt{a - \sqrt{a + x}} = x$$
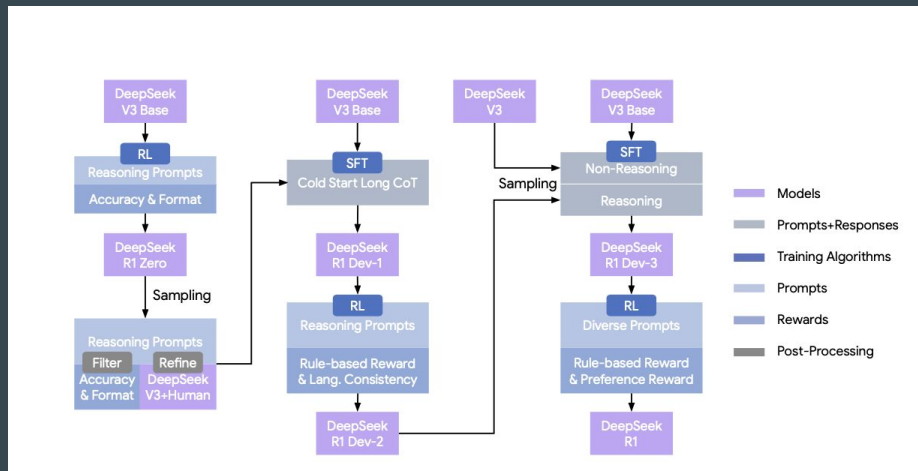
First, let's square both sides:

$$a - \sqrt{a + x} = x^2 \implies \sqrt{a + x} = a - x^2$$

Next, I could square both sides again, treating the equation: $\cdots$
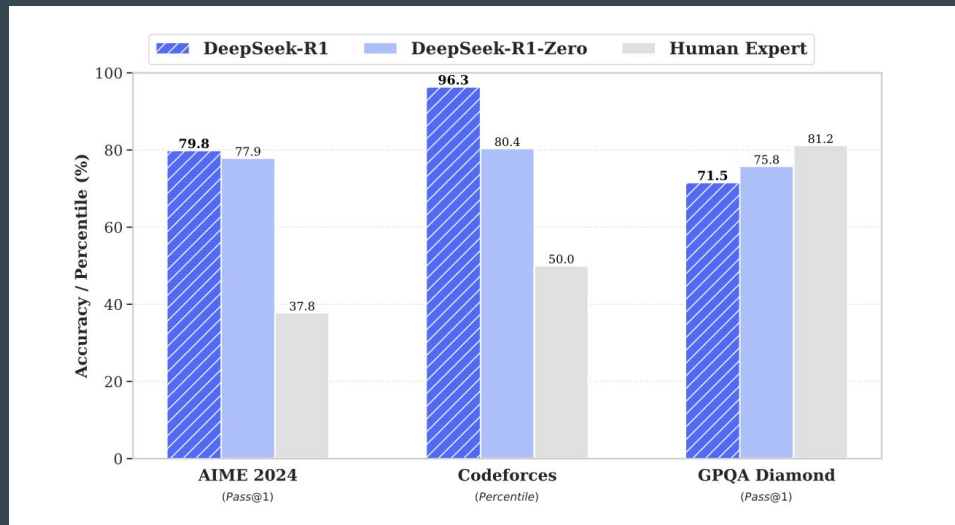
$\cdots$

# The DeepSeek-R1 Training Pipeline

- Stage 1 (Cold Start): SFT on ~5k long-COT samples for readability

- Stage 2 (Reasoning RL): Large scale RL focused on STEM + Logic

- Stage 3 (Rejection Sampling): Synthetic data generation (800k samples)

- Stage 4 (Preference RL): Alignment for helpfulness and safety

# Benchmarking Reasoning Capabilities

- Competitive Math: 79.8% on AIME 2024

- Competitive Coding: 96.3% human percentile on Codeforces

- Conclusion: Achieves similar performance with OpenAI-o1-1217 using significantly less compute

# Knowledge Distillation

- Data Gathering: Used DeepSeek-R1 to generate ~800k high-quality "reasoning traces"

- Filtering: Only "correct" answers with readable logic were kept

- Fine-Tuning: These traces were used for Supervised Fine-Tuning on smaller models

- Efficiency: A distilled 1.5B model outperforms non-reasoning GPT-4o on mathematical benchmarks

# Safety Evaluation and Risk Control

- Point: The model implements a "Safety Reward" during the final RL stage to mitigate risks without compromising reasoning

- System: Use of an external risk review system (sending queries to DeepSeek-V3) to detect and retract unsafe dialogue
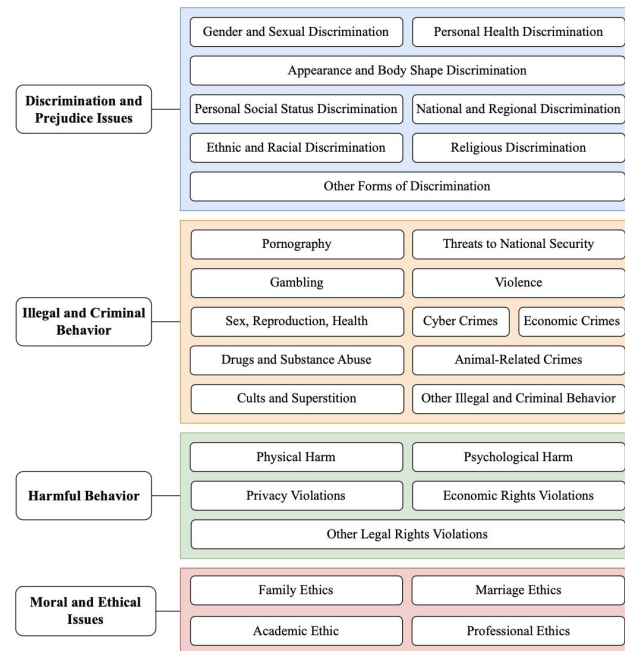


Figure 13 | Taxonomy of in-house safety benchmark.

# Unsuccessful Approaches

- Process Reward Models (PRM): Deemed inefficient due to the difficulty of defining universal "correct steps" and the high risk of reward hacking

- Monte Carlo Tree Search (MCTS): Challenges in scaling due to the exponentially larger search space of token generation

# Limitations

- Inference Efficiency: Models lack an internal "stop signal," often overthinking trivial tasks and wasting compute

- Domain Constraints: Current success is tied to "verifiable" tasks (STEM) and extending this to subjective fields remains a challenge

- Linguistic Consistency: Occasional language mixing and drops in consistency occur during long, multi-turn interactions

- Prompt Sensitivity: Performance can be inconsistent based on phrasing, and the model often requires explicit instructions to "think" for optimal results

# Future work

- Adaptive Compute: Scaling reasoning time based on task difficulty to improve efficiency

- Agentic Tasks: Expanding from logic puzzles to real world engineering and autonomous planning

- Multi-Modal Reasoning: Integrating reasoning chains with vision and audio inputs for cross-modal tasks

- Nuanced Rewards: Developing reward systems for subjective tasks that lack binary answers

# Conclusion

- Emergent Reasoning: Proves logic (backtracking, self-correction) can be discovered via RL without human data
- Efficiency: GRPO eliminates the Value Model, drastically reducing the hardware cost of training
- Accessible AI: Distillation successfully transfers traces to small, efficient models
- A New Scaling Rule: Intelligence now scales with test-time compute rather than just model size.

DeepSeek-R1 provides a blueprint for next-gen AI, replacing human-led imitation with autonomous, machine native logic

# Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model?
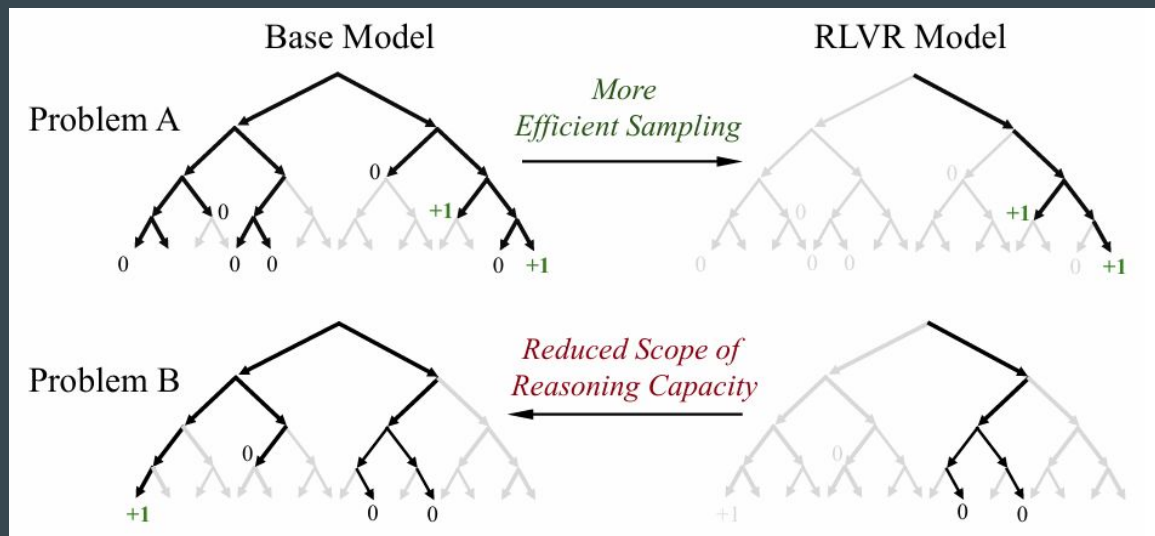
Yue et al.
arXiv:2504.13837

# Overview

- Common belief: RL helps agents discover new strategies
  - Go
  - Atari

- People expect RL with verifiable rewards (RLVR) to "self-improve" like classic RL

Main Research Question: Does current RLVR discover *new* reasoning patterns, or just exploit base-model ones?

# Core Intuition: Efficiency vs Boundary

- Sampling efficiency: how likely the model reaches a correct path in few tries
- Reasoning boundary/coverage: what fraction of problems are solvable by the model given many tries

# Key Metric: pass@k

pass@k = "can we solve it if we try enough times?"

- Sample k outputs; pass@k = 1 if any of them are correct
- Estimated by sampling k completions using temperature-based sampling

Purpose:

- Evaluating boundaries
- If RL truly adds reasoning ability, it should solve problems the base model can't
- **Efficiency** = better pass@1
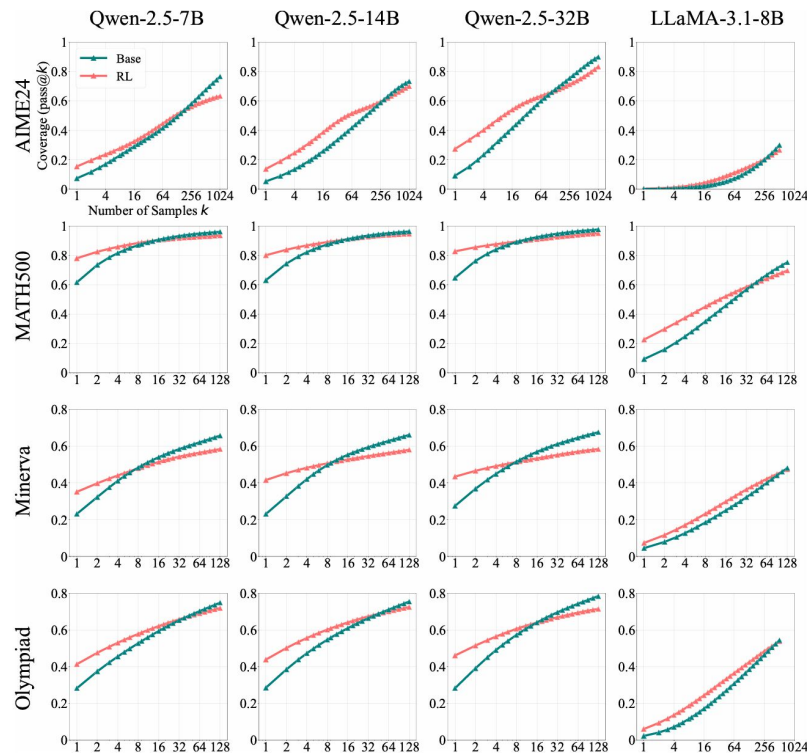- **Boundary** = better pass@k as k → large

# Experiments at a Glance

- Test across multiple model families + RL algorithms
- Goal: find general patterns across several types

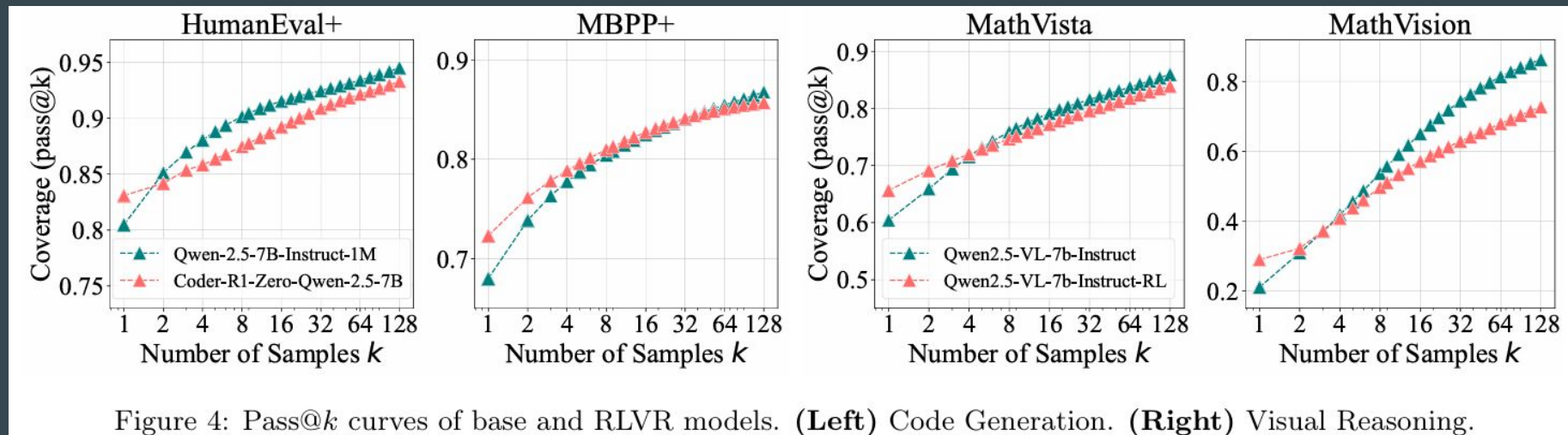| Task | Start Model | RL Framework | RL Algorithm(s) | Benchmark(s) |
|------|-------------|--------------|-----------------|--------------|
| **Mathematics** | LLaMA-3.1-8B<br>Qwen2.5-7B/14B/32B-Base<br>Qwen2.5-Math-7B | SimpleRLZoo<br>Oat-Zero<br>DAPO | GRPO | GSM8K, MATH500<br>Minerva, Olympiad<br>AIME24, AMC23 |
| **Code Generation** | Qwen2.5-7B-Instruct<br>DeepSeek-R1-Distill-Qwen-14B | Code-R1<br>DeepCoder | GRPO | LiveCodeBench<br>HumanEval+ |
| **Visual Reasoning** | Qwen2.5-VL-7B | EasyR1 | GRPO | MathVista<br>MathVision |
| **Deep Analysis** | Qwen2.5-7B-Base<br>Qwen2.5-7B-Instruct<br>DeepSeek-R1-Distill-Qwen-7B | VeRL | PPO, GRPO<br>Reinforce++<br>RLOO, ReMax, DAPO | Omni-Math-Rule<br>MATH500 |

# Main result 1: Math pass@k curves flip

- RLVR > base at small k (pass@1)
- Base > RLVR at larger k

Bottom Line: RL helps get answers faster, but the base model often has *broader coverage* when you sample enough

# Main result 2: Same pattern in code + visual



Figure 4: Pass@k curves of base and RLVR models. **(Left)** Code Generation. **(Right)** Visual Reasoning.

Same behavior found outside math. Base model surpasses RLVR at higher k, suggesting RLVR isn't expanding set of solvable problems

# Then why do accuracy averages increase?

Observations:
- More mass near accuracy ~1.0
- Fewer low accuracies
- BUT more accuracy = 0

Interpretation:
- RLVR makes some problems unsolvable
- Polarizes probabilities; hyper-focuses on already good strategies



Figure 5: Qwen2.5-7B Accuracy Histogram on Minerva.

# Coverage analysis: RLVR solves (almost) a subset of base model

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

Very few cases where RLVR solves something the base model can't, but many the other way around.

# Perplexity test: Are RL solutions already in the base model?

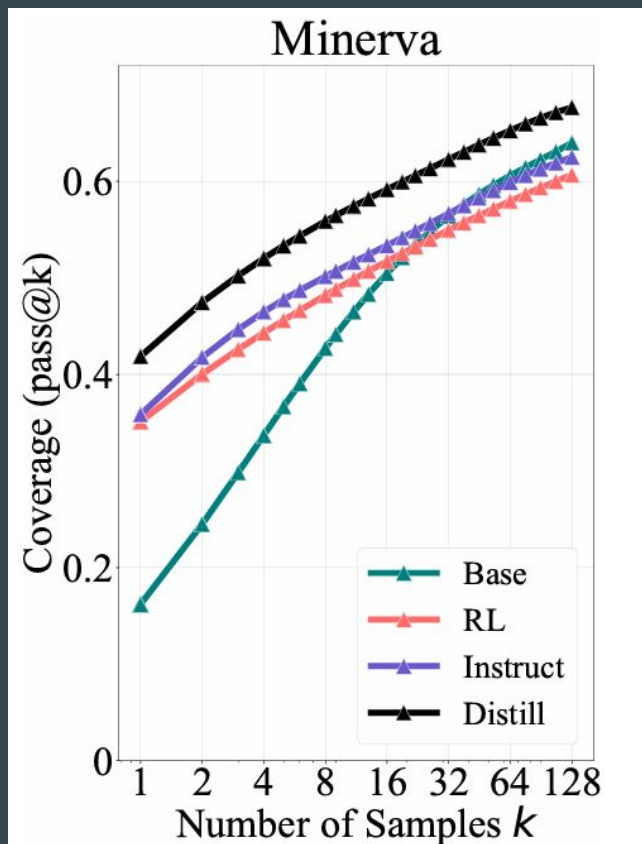Perplexity: if base assigns low perplexity to RL-generated solutions, they are likely to be generated by the base

Observation: Base model is not surprised by RL's answers, so RL is exploiting existing reasoning paths

$$\mathrm{PPL}_m(\mathbf{Y} \mid x) = \exp\left(-\frac{1}{T}\sum_{t=1}^{T}\log P(y_t \mid x, y_1, \ldots, y_{t-1})\right)$$

# Distillation is different: it *can* expand the boundary

- Distilled model pass@k consistently above base model
- Imports new reasoning patterns from a stronger teacher
- Truly expands the student's reasoning boundary



Minerva

Coverage (pass@k) vs Number of Samples $k$, with curves for Base, RL, Instruct, and Distill.

# Discussion: why RLVR gets stuck + how to fix it

1. Vast action space
   a. Exponentially larger than other environments
   b. Effective exploration is much harder
2. Starting from a pre trained prior:
   a. Can't start from scratch like other RL environments
   b. Struggles to explore new reasoning patterns
   c. Samples that deviate almost always get negative rewards

Future Directions:

- Fine-grained credit assignment – add intermediate signals, not just 0/1
- Scaling RL iteration – start by training on easier subproblems

# Conclusion and Limitations

Main Ideas:

1. RLVR boosts pass@1 efficiency
2. Base can surpass at large k → boundary not expanded
3. RL reasoning paths often already present in base model
4. Re-weights distribution making some reasoning paths extremely unlikely

Limitations:

1. Many top models and training pipelines are proprietary
2. Field is moving fast, emerging techniques may solve this

# ProRL: Prolonged Reinforcement Learning Expands Reasoning Boundaries in Large Language Models

Liu et al.
arXiv:2505.24864

# Background

- ProRL Methodology: Is there still any merit to RLVR training?

- Updated Training Scheme

- What was Accomplished: Can models learn something new through the use of RLVR?

# Methodology: Mitigating Entropy Collapse

$$\text{clip}(r_\theta(\tau), 1 - \epsilon_{low}, 1 + \epsilon_{high}).$$

- **Decoupled Clipping**
  - Separates the Upper and Lower Clipping Bounds (Token Probability)
  - What does this accomplish?
- **Dynamic Sampling Policy Optimization**
  - Filtering Questions Based on Performance
- **Reference Policy Reset**
  - KL Vanishing Problem
  - KL Penalty
  - Resetting the "Anchor Point"
- **Prolonged Training**



Figure 2: ProRL training dynamics.

$$L_{KL-RL}(\theta) = L_{GRPO}(\theta) - \beta D_{KL}(\pi_\theta || \pi_{ref}).$$

# Task Diversity

- **Math (AIME/AMC):**  Provides strict, binary rewards to teach precision.

- **Code (Codeforces)**  : Provides continuous rewards (passing test cases) to teach debugging and edge-case handling.

- **Logic (Reasoning Gym)**  : Targets pure reasoning rules (e.g., Zebra puzzles) where base models are weakest.

- **STEM**: Tests Knowledge Application. Combines Specific Scientific Knowledge with reasoning for better real-world problem representation

- **Instruction Following:**  Can the model continue to follow basic instruction

# Experimental Setup

- Baseline & Competitors:

  - **The Base: DeepSeek-R1-Distill-Qwen-1.5B** (The starting student model).

  - **The Specialists: DeepScaleR (Math-only) and DeepCoder (Code-only)** of the same size. ProRL beats them both (+4.6% Math, +6.5% Code).

  - **The Giant: DeepSeek-R1-Distill-Qwen-7B (5x larger)**. ProRL matches or beats this larger model despite having fewer parameters.

# Key Results

- Reversing Collapse
  - Widening Distribution of "Correct Answers"
- Weak-to-Strong Generalization
  - Where does RL work best? Where is it less effective?
  - Zero-to-One: OOD capability achieved
- Creativity Index
  - DOLMA: pre training corpus
  - N-gram Overlap
  - Where was their little Creativity?




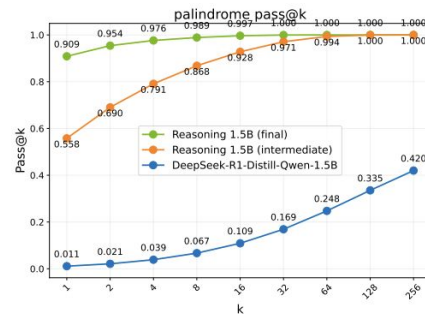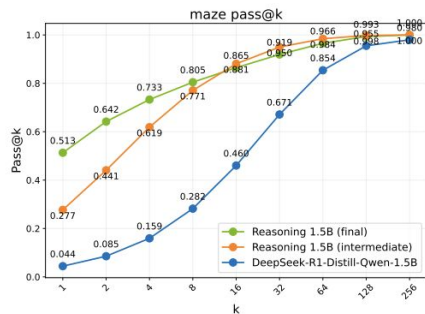
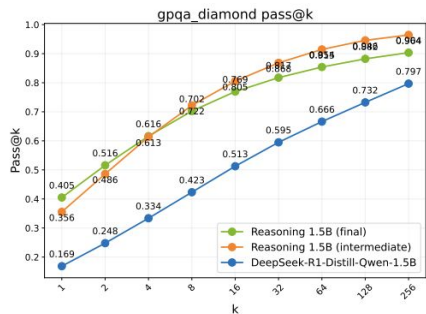Figure 6: ProRL generalizes to increased task difficulty on *graph_color*.



Figure 5: Expanded reasoning boundary for OOD task *boxnet*.

Mean: 3.84    Mean: 4.42    Mean: 4.70
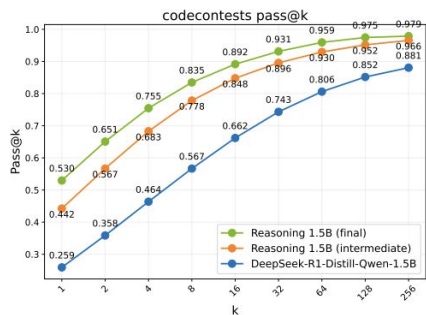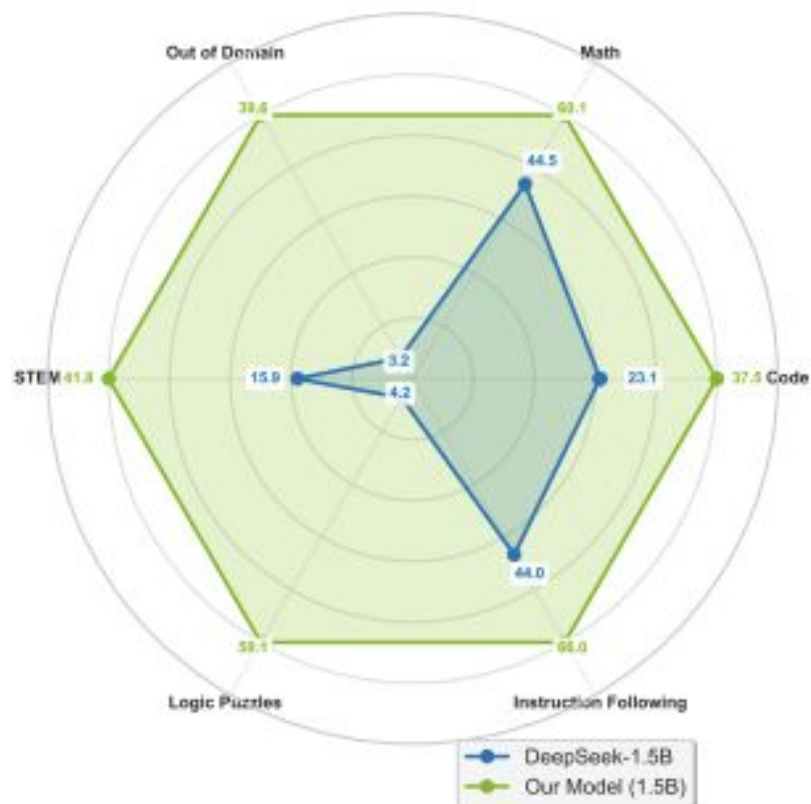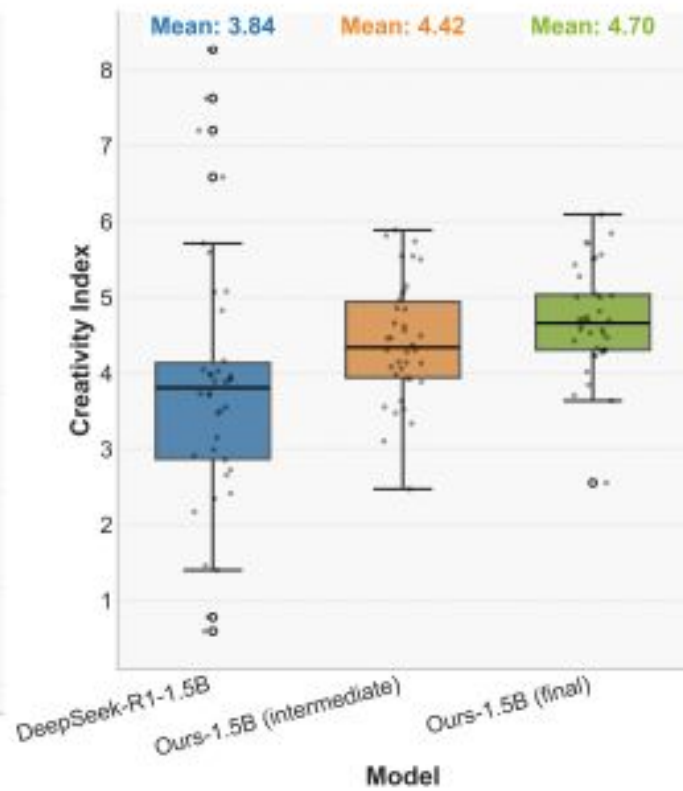
# Conclusion

- **Conclusive Results** : Stable and prolonged RL training expands reasoning capabilities.

- **Effective RL Training Requires** : Long training, reference resets, and diverse tasks.

- **Future Impact** : Small models can achieve powerful reasoning, democratizing high-level AI access.

# Limitations

- **Limited Task Scope:** Small subset of all tasks

- **Scalability Concerns:** Work done on smaller models

- **Computational Resources:** Extensive Training = Great Computational Demands

- **Unstable Training:** Hard Resets leads to variability in repeated trials

# Culminating Conclusion

- **RLVR works:**  big gains on verifiable reasoning tasks (math/code)
- But "reasoning gains" could either mean efficiency or capability
- **DeepSeek-R1:**  RLVR can elicit longer, self-correcting reasoning behavior
- **Boundary paper:**  boosts pass@1 but may not expand and can even shrink reasoning boundaries
- **ProRL:** prolonged/stabilized RL can expand reasoning boundaries

**Main Takeaway:**  RLVR is rapidly evolving; results depend on the training pipeline, and new research is constantly improving outcomes

# Thank you. Any Questions?