

MATH-BEYOND: A BENCHMARK FOR RL TO EXPAND BEYOND THE BASE MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

With the advent of DeepSeek-R1, a new wave of reinforcement learning (RL) methods has emerged that seem to unlock stronger mathematical reasoning. However, a closer look at the open-source ecosystem reveals a critical limitation: with sufficiently many draws (e.g., $\text{pass}@1024$), many existing base models already solve nearly all questions on widely used math benchmarks such as MATH-500 and AIME 2024. This suggests that the RL fine-tuning methods prevalent in the LLM reasoning literature largely sharpen existing solution modes rather than discovering entirely new ones. Such sharpening stands in contrast to the broader promise of RL: to foster exploration and to acquire new skills. To move beyond this plateau, we introduce MATH-Beyond (MATH-B), a benchmark deliberately constructed to defeat common open-source models of up to 8B parameters even under large sampling budgets. Improving performance on our benchmark via RL requires methods that learn to reason in ways that go beyond base model capabilities in repeated sampling. Since the problems are drawn from subsets of DAPO-Math-17K and DeepScaleR datasets, they remain topically equivalent to standard high-school math. Validating our premise, RL fine-tuned models such as Nemotron-Research-Reasoning-Qwen-1.5B and DeepScaleR-1.5B-Preview perform poorly on MATH-B at $\text{pass}@1024$, showing how existing approaches fall short on tackling harder instances. We hope MATH-B will catalyze exploration-driven RL approaches that elicit deeper reasoning capabilities. We release MATH-B at <https://huggingface.co/datasets/brendel-group/MATH-Beyond>.

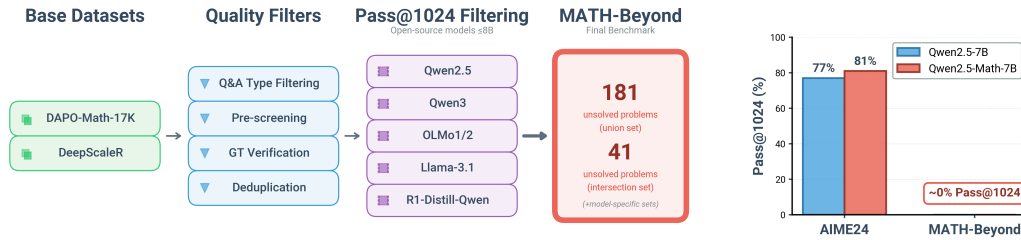


Figure 1: **MATH-Beyond: Benchmark Construction and Difficulty.** *Left:* Schematic of the MATH-B creation process. A large set of problems from DAPO-Math-17K and DeepScaleR is first refined through quality filters to ensure answer correctness and verifiability. This is followed by evaluation against a gauntlet of open-source base models ($\leq 8B$, e.g., Qwen3, Qwen2.5 (-Math), DeepSeek-R1-Distill), at a $\text{pass}@1024$ budget to isolate problems that lie beyond their limits. The filtering yields the MATH-B suite of benchmarks: a 41-problem intersection set (unsolved by all base models) for evaluating universal difficulty, and a larger 181-problem union set (unsolved by at least one model) with model-specific splits for targeted analysis. This suite provides a rigorous testbed to drive the development of exploration methods for RL. *Right:* An illustration of MATH-B’s significant difficulty compared to common test sets like AIME24. Representative open-source models like Qwen2.5 achieve near-zero $\text{pass}@1024$ scores on MATH-B, highlighting its difficulty. Qwen2.5 results are from Yue et al. (2025).

1 INTRODUCTION

In the 2010s, deep reinforcement learning showcased its power through striking demonstrations of exploration and skill acquisition (Mnih et al., 2013). Atari agents, starting from random play, mastered complex games by discovering strategies unreachable to base policies, guided by exploration incentives and intrinsic rewards (Ladosz et al., 2022; Amin et al., 2021). Methods such as count-based exploration (Bellemare et al., 2016) and later Go-Explore (Ecoffet et al., 2021) drove dramatic jumps from inept play to expertise, highlighting RL’s ability to uncover new capabilities. Around the same time, AlphaGo (Silver et al., 2016) and AlphaZero (Silver et al., 2017) extended this promise to board games like Go, Chess, and Shogi, where self-play took agents from scratch to superhuman mastery, revealing novel strategies along the way.

Against this backdrop, academic progress in reasoning-focused LLMs has taken a very different path. Community-trained models often show improved accuracy on popular benchmarks such as MATH or AIME24 (Liu et al., 2025b; Song et al., 2025; Chen et al., 2025; Cheng et al., 2025; Cui et al., 2025; Shao et al., 2025; Wang et al., 2025; Yu et al., 2025b). However, these RL models typically succeed only on problems that their corresponding base models could already solve given realistic sampling budgets (Wu et al., 2025; Yue et al., 2025). This is a far cry from earlier RL successes, where base policies were incapable of solving tasks outright and progress required genuine exploration and skill acquisition. This disconnect between RL’s exploratory promise and its current application reflects a substantial blindspot in the current open-source evaluation ecosystem. Because several open-source base models already achieve nearly 100% `pass@1024` on several popular benchmarks (Yue et al., 2025), test sets in their current form are fundamentally inadequate for measuring—or encouraging—genuine progress in reasoning beyond the base model’s reach.

To address this gap, we introduce **MATH-Beyond (MATH-B)**, a new benchmark of high-school-level competition math problems, specifically constructed so that popular open-weight base models are unlikely to solve even with 1024 attempts. As a result, progress on MATH-B necessarily requires expanding the reasoning capabilities of base models, making it an ideal target for academic research.

MATH-B is constructed by filtering mathematical reasoning datasets (DAPO-Math-17K (Yu et al., 2025b) and DeepScaleR (Luo et al., 2025)), resulting in problems that are topically indistinguishable from those in standard benchmarks. While constructing the dataset, we also uncovered and addressed several non-obvious failure modes in programmatic verification, which informed our final benchmark design. To ensure correctness, all problems are additionally verified against stronger reasoning models such as GPT-5-Mini or o4-mini-high, which reliably solve them. We further confirm that leading community RL models, including Nemotron-Research-Reasoning-Qwen-1.5B (Liu et al., 2025a) and DeepScaleR-1.5B, perform poorly on MATH-B, underscoring the limitations of current approaches and the need for methods that extend reasoning capabilities. In summary, our contributions are as follows:

- **New benchmark suite:** We construct MATH-Beyond, a benchmark suite derived from the failures of a large and diverse set of base models (`pass@1024` ≈ 0). This suite includes: a union set of 181 problems unsolved by at least one of these models; model-specific benchmarks for targeted analysis; and a highly challenging intersection set of 41 problems that proved unsolvable for the entire considered set. To ensure quality, all problems are annotated for topic and difficulty following the procedure from Omni-MATH (Gao et al., 2024a), and undergo answer verification by frontier models (o4-mini-high & GPT5-Mini).
- **Verification pitfalls:** During benchmark construction, we observed several pitfalls in standard RLVR verification. We take these into account in our benchmark design and also document them, highlighting subtle edge cases for the community to be aware of.
- **Model evaluation:** We evaluate RL-finetuned models such as Nemotron-Research-Reasoning-Qwen-1.5B, DeepScaleR-1.5B, and Skywork-OR1 (He et al., 2025a) on MATH-B and find that they do not substantially expand reasoning boundaries. In contrast, we find that newer model families like Qwen3-4B and Qwen3-8B (Yang et al., 2025) perform better, presumably due to better distributional overlap with our dataset.

2 FRAMEWORK FOR EVALUATING MODEL EXPANSION

Our goal is to quantify whether a post-trained model (e.g., after RLVR) *expands* its reasoning capabilities beyond its base model. While post-training can affect many aspects (e.g., robustness, exploration, or transfer), our framework is designed to isolate the specific phenomenon of *boundary expansion*—what new problems the post-trained model can solve relative to its base model. This section formalizes this evaluation framework, adapting the nomenclature and definitions from Wu et al. (2025) into a slightly simplified, empirical version. We apply this framework to our benchmark MATH-B, which is a “zero-baseline” test. We define this as a setting where a benchmark is specifically constructed such that the base model empirically fails on all problems (i.e., its observed pass@k is zero) given a realistic sampling budget.

2.1 THE EMPIRICAL PASS@K METRIC

We evaluate a post-trained policy π against its base model q on a dataset D . Our evaluation is based on the empirical pass@k metric. For a given problem $x \in D$ and a policy $p \in \{\pi, q\}$ (usually an LLM), we draw k i.i.d. samples $\{y_1, \dots, y_k\}$. Let $\mathcal{C}(x)$ be the set of all correct completions for x . The empirical success is:

$$\text{pass@k}(p; x) = \begin{cases} 1 & \text{if } \exists i \in \{1, \dots, k\} \text{ such that } y_i \in \mathcal{C}(x), \\ 0 & \text{otherwise.} \end{cases}$$

This binary value—1 for a “pass” and 0 for a “fail”—is the ground-truth measure of success for a single problem.

With a slight abuse of notation, the central metric reported in our paper is the average success rate across the entire dataset D :

$$\text{pass@k}(p) = \frac{1}{|D|} \sum_{x \in D} \text{pass@k}(p; x),$$

the **ratio of problems in D solved by policy p** using a k -sample budget.

2.2 DECOMPOSITION AND KEY METRICS

To understand *how* π ’s performance differs from q ’s, we first define the **Reachable Set** $\mathcal{R}_k(p, D)$ as the set of problems p solves:

$$\mathcal{R}_k(p, D) = \{x \in D : \text{pass@k}(p; x) = 1\}.$$

Comparing the reachable sets of π and q allows us to isolate our primary metric and contextualize it.

Expansion (Primary Metric). Our *primary focus* and sole reported metric is the **Expansion Rate**. This measures genuine boundary expansion—new problems π solves that q could not. It is defined based on the **Expansion Set** $\mathcal{E}_k = \mathcal{R}_k(\pi, D) \setminus \mathcal{R}_k(q, D)$:

$$\text{Expansion Rate} = \frac{|\mathcal{E}_k|}{|D|}.$$

Shrinkage. For diagnostic context, we also define *Shrinkage* (or “forgetting”). This measures problems q could solve that π now fails, defined by the *Shrinkage Set* $\mathcal{S}_k = \mathcal{R}_k(q, D) \setminus \mathcal{R}_k(\pi, D)$. It can be quantified as:

$$\text{Shrinkage Rate} = \frac{|\mathcal{S}_k|}{|D|}.$$

Preservation. Similarly, *Preservation* measures the fraction of q ’s capabilities that π retains, defined by the *Preservation Set* $\mathcal{P}_k = \mathcal{R}_k(\pi, D) \cap \mathcal{R}_k(q, D)$. It is quantified as:

$$\text{Preservation Rate} = \frac{|\mathcal{P}_k|}{|\mathcal{R}_k(q, D)|}.$$

Consolidation. Finally, *Consolidation* is a concept for measuring if preserved solutions become more robust (i.e., solvable at `pass@1`). It is defined as:

$$C_k(\pi, q) = \frac{|\mathcal{P}_k \cap \mathcal{R}_1(\pi, D)|}{|\mathcal{P}_k|}.$$

Interpretation. The overall pass rate of π decomposes as $\text{pass}@k(\pi) = (|\mathcal{E}_k| + |\mathcal{P}_k|)/|D|$. While Shrinkage, Preservation, and Consolidation are crucial concepts for a full diagnosis, our target for *expanding the reasoning boundary* is squarely the **Expansion Rate**. The other concepts serve as a theoretical guardrail to ensure that measured gains are from genuine expansion, not mere reallocation.

2.3 SPECIAL CASE: THE MATH-B BENCHMARK

The metrics framework in Section 2.2 simplifies cleanly for MATH-B. This benchmark is constructed as a “zero-baseline” meaning it is composed of problems where the base model q was empirically observed to fail within the sampling budget k (i.e. $\text{pass}@k \approx 0$). In our evaluation, this means the base model’s reachable set is empty::

$$\mathcal{R}_k(q, D) = \emptyset.$$

Under this premise, the decomposition from Section 2.2 collapses:

$$\mathcal{S}_k = \emptyset, \quad \mathcal{P}_k = \emptyset, \quad \mathcal{E}_k = \mathcal{R}_k(\pi, D).$$

Shrinkage and preservation are thus not applicable. Every solve by π is by definition an expansion, and the single relevant metric becomes:

$$\text{Expansion Rate} = \frac{|\mathcal{R}_k(\pi, D)|}{|D|} = \text{pass}@k(\pi).$$

A “win” on MATH-B is simply a positive Expansion Rate, providing an unambiguous readout of genuine boundary expansion.

3 CREATING THE BENCHMARK

3.1 ON THE PITFALLS OF VERIFICATION

The goal of our work is to identify questions that models fail at high `pass@1024` due to genuine methodological or calculation errors, not because of artifacts in the verification pipeline. Current rule-based verification methods for math, used both in training and evaluation, are themselves prone to systematic failures. These failures can mask true performance by penalizing answers where the reasoning or ground truth is technically correct but formatted poorly, making it appear as though the model failed when the verifier simply could not handle the formatting or parsing.

To guide our evaluation and benchmark design, and to alert the community to these pitfalls, we sketch seven distinct failure modes observed across common RL finetuning and evaluation frameworks. We uncovered these issues through an analysis of reasoning traces produced by DeepSeek-R1-Distill-Qwen2.5-7B (DeepSeek-AI et al., 2025) when applied to a subset of NuminaMath-1.5 (LI et al., 2024). Table 1 provides concise but complete examples (problem, ground truth, model output snippet, and the precise mismatch), and Table 3 in the Appendix summarizes how different frameworks are affected.

3.2 MATH-B CONSTRUCTION PIPELINE

Our benchmark is the result of a multi-stage filtering pipeline designed to isolate problems that are (a) correct and unambiguous, (b) novel, and (c) demonstrably unsolvable by a suite of open-source base models even at a high sampling budget. The process is detailed below and in Figure 1.

3.2.1 BASE DATASET SOURCING

Our construction begins by sourcing a large pool of candidate problems. The primary goal, as stated, is to find problems that defeat open-source models at high sampling budgets (`pass@1024`). This

Table 1: **Common failure modes in rule-based math answer verification.** These failures often stem from rigid heuristics, such as reading only the first or last boxed answer, requiring specific text anchors (e.g., "Answer:"), or other parsing failures. Each row shows the ground truth (GT), a model snippet, and the resulting verifier error.

Mode	GT	Model snippet	Verifier behavior
F1: Multiple valid	5 or 13	... $x = \boxed{5}$ or $\boxed{13}$	<i>Expected:</i> accept both values. <i>Actual:</i> only first/last boxed read; one valid ignored.
F2: Late correct	150	... $\boxed{50}$; ... $\boxed{100}$; ... $\boxed{150}$	<i>Expected:</i> final value only. <i>Actual:</i> early intermediate captured.
F3: Early correct	42	... $\boxed{42}$... later $\boxed{84}$	<i>Expected:</i> accept any correct boxed answer. <i>Actual:</i> later wrong overrides early correct.
F4: Corrections	25	... $\boxed{20}$, then correcting $\boxed{25}$	<i>Expected:</i> use corrected value. <i>Actual:</i> first match persists; "last-only" works, "first-only" fails.
F5: Unordered tuple	3, 4, 5	... $\boxed{5}$, $\boxed{3}$, $\boxed{4}$	<i>Expected:</i> order-agnostic match. <i>Actual:</i> order-sensitive; heuristics misgrade.
F6: Missing anchors	4	... $2 + 2 = \boxed{4}$ (no "Answer:" prefix)	<i>Expected:</i> accept boxed numeral alone. <i>Actual:</i> anchor required; otherwise wrong.
F7: MCQ partial	C) 1000	... "The answer is \boxed{C} ."	<i>Expected:</i> accept label-only (C). <i>Actual:</i> requires exact "C)1000".

requirement immediately disqualifies many common datasets. For instance, DeepMath103K (He et al., 2025b) and Big-Math (Albalak et al., 2025) are explicitly designed to be solvable (e.g., for GRPO training), and indeed, models like Llama-3.1-8B solve them with few attempts. We also sought to minimize data contamination, ruling out datasets like NuminaMath which have been seen in pretraining by models like Qwen2.5-Math (Yang et al., 2024) and consequently DeepSeek-R1.

We therefore selected two base datasets: **DAPO-Math-17K** (Yu et al., 2025a), which satisfies our criteria for both high difficulty (not verified by open-source models) and novelty (released after DeepSeek-R1), and **DeepScaleR** (Luo et al., 2025), which, while potentially seen, provides a large corpus of problems that also lack open-source verification. This combined set forms our initial candidate pool of 53,682 problems.

3.2.2 QUALITY FILTERS

We then apply an array of quality filters as described in the following paragraphs.

Question and Answer-Type Filtering To mitigate the verification pitfalls detailed in Section 3.1, we apply several deterministic filters. First, to avoid ambiguity from failure modes like F1 (Multiple valid answers), and F5 (Unordered Tuples), we filter DeepScaleR to include only problems with integer-based ground truths (DAPO-Math-17K already meets this criterion). Second, using regex-based filters, we remove all multiple-choice questions (to prevent F7 parsing failures) and any questions containing Chinese characters. Finally, to ensure all problems are self-contained, we remove any questions referencing external figures or images. This filtering process reduced the pool to 34,515 datapoints.

Pre-screening and Random Sampling The problems from the previous step still represent a computationally prohibitively large set for large-scale evaluation. To reduce this pool, we first perform a difficulty pre-screening step. We evaluate Deepseek-R1-Distill-Qwen2.5-7B on all problems,

keeping only those that it *could not solve* within a `pass@16` budget. In addition to the screening, we randomly sample a portion of this dataset for further processing.

Ground-Truth Answer Verification A critical step is to ensure that these problems are unsolved due to their intrinsic challenge, not because their provided ground-truth answers are incorrect. While our source datasets are generally reliable (e.g., DeepScaleR is derived from AIME (MAA) and AMC (MAA, 2023) competitions (Luo et al., 2025)), we took an additional step to ensure the ground-truth answers are correct. We evaluate (`pass@2`) this pre-screened subset using two frontier-class models, o4-mini-high and GPT-5-Mini (Balunović et al., 2025). For each problem, we prompt the models to *"Think step-by-step and put your final answer in \boxed."* and check if the extracted answer matches the dataset’s ground truth. We retain only the problems where *at least one* of these frontier models successfully reproduced the ground-truth answer.

Deduplication Against Standard Benchmarks To ensure the novelty of our dataset, we first perform an exact string-match deduplication of our candidate problems against several common test sets, including MATH-500 (Hendrycks et al., 2021), MinervaMath (Lewkowycz et al., 2022), OlympiadBench (He et al., 2024), AMC23 (MAA, 2023), AIME-2024, and AIME-2025 (MAA). We confirmed that no question from our set is present in these benchmarks. At the end of this stage, a collection of 184 problems remain.

3.2.3 FINAL BENCHMARK CONSTRUCTION (`PASS@1024` FILTERING)

The final step in our benchmark’s construction is to filter the candidate problems using a diverse suite of representative open-source models. The definition of a *‘base model’* is contextual; it typically refers to a model intended for further fine-tuning. Keeping in mind the models the community often uses for post-training research, we select a suite of models categorized as either ‘base’ or ‘supplementary’ models. As we will detail, these two groups are used to construct different subsets of our final benchmark.

Base Models This set is used to define the most challenging subset of our benchmark. It includes: Qwen2.5-1.5B, Qwen2.5-7B (Qwen et al., 2025), Qwen2.5-Math-1.5B, Qwen2.5-Math-7B (Yang et al., 2024), Qwen3-4B-Base, Qwen3-8B-Base (Yang et al., 2025), DeepSeek-R1-Distill-Qwen2.5-1.5B, DeepSeek-R1-Distill-Qwen2.5-7B (DeepSeek-AI et al., 2025), OLMo-7B (Groeneveld et al., 2024), OLMo-2-7B (OLMo et al., 2025), and Llama-3.1-8B (Grattafiori et al., 2024).

Supplementary Models This group is combined with the base models to define the full benchmark. It includes: Qwen2.5-1.5B-Instruct, Qwen2.5-7B-Instruct (Qwen et al., 2025), Qwen2.5-Math-1.5B-Instruct, Qwen2.5-Math-7B-Instruct (Yang et al., 2024), Qwen3-4B, Qwen3-8B (Yang et al., 2025), DeepScaler-1.5B (Luo et al., 2025), Nemotron-Research-Reasoning-Qwen-1.5B(v1 and v2) (Liu et al., 2025a), and Skywork-OR1-7B (He et al., 2025a).

The deduplicated and pre-screened candidate set is then subjected to our final filtering stage. We evaluate the problems against *all listed models* generating 1024 samples for each. During this evaluation, we apply our robust verification logic (see Table 1 for failure modes) to ensure we were measuring genuine reasoning failures. Our experiments and analyses required over 20 000 A100 GPU hours. Please refer to Appendix A.2 for more details on the evaluation parameters.

This comprehensive evaluation allows us to define our final benchmark, **MATH-B**, which comprises three broad subsets for targeted analysis:

- **MATH-Beyond-Union Set (MATH-B-U):** The full benchmark of **181 problems**, containing any problem that at least one model from our *entire suite* (both base and post-trained) failed to solve within 1024 attempts.
- **MATH-Beyond-Intersection Set (MATH-B-I):** A more challenging core subset of **41 problems**, that *all* of our considered **base models** failed to solve. This is a hard subset of MATH-B-U (see Table 5 for the QA pairs).
- **Model-Specific Sets:** For any given model in our suite, this is the collection of all problems in the Union Set that it failed to solve (see Table 4). These sets enable fine-grained testing of RL finetuned models derived from one of the evaluated base models.

Overall, our MATH-B datasets serve as benchmarks of reasoning capability expansion (see Section 2).

3.3 DATA CHARACTERISTICS AND KEY REMARKS

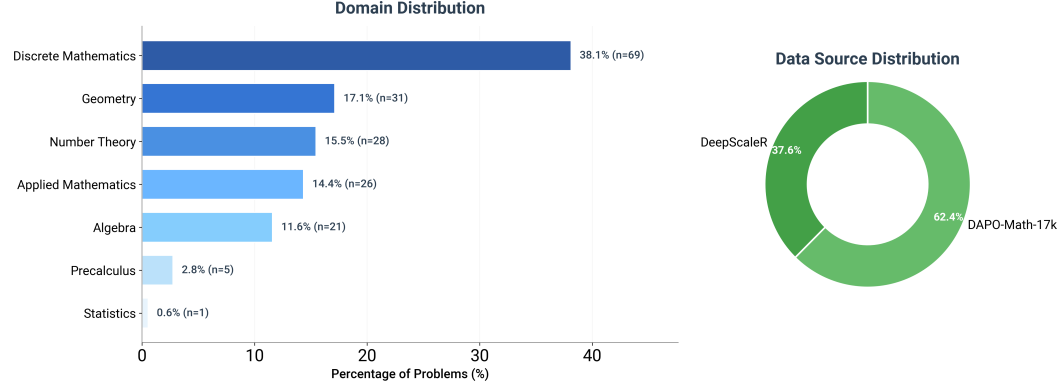


Figure 2: **Characteristics of the MATH-B-U dataset.** Left subplot shows the distribution of math domains. Right subplot show the distribution of source datasets.

Domain and Difficulty Annotation To analyze the characteristics of the MATH-B dataset, we annotate each problem for its mathematical domain and human-perceived difficulty. We adapt the procedure from Omni-MATH (Gao et al., 2024a), using GPT-5 to perform the labeling based on a contrastive prompting strategy (see supplementary material for the full prompt). This method leverages labeled examples from various math datasets to contextualize and assign scores to new problems and has shown to be aligned well with human judgment (Gao et al., 2024a). A selection of annotated problems from MATH-B-I is provided in Table 5. We also plot the distribution of ground truth answers in Figure 4.

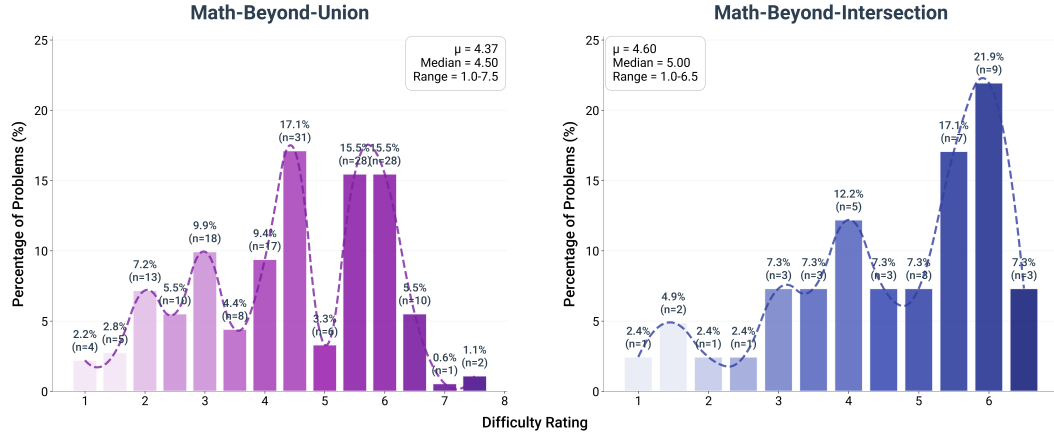


Figure 3: **Difficulty distribution.** The left subplot shows the difficulty distribution of the MATH-Beyond-Union set, while the right subplot shows that of the MATH-Beyond-Intersection set. The wide spread of difficulty levels highlights a key mismatch: the problems that models find challenging are not necessarily those that humans typically struggle with.

Distribution Analysis As shown in Figure 2, the topic distribution of MATH-B consists entirely of standard high-school mathematics subjects, ensuring topical relevance. The difficulty distribution, plotted in Figure 3, reveals a wide range for both the Union and Intersection sets, with a median human-difficulty rating of 4 out of 10. Notably, even for the challenging MATH-B Intersection set, the maximum difficulty score is only 6.5. This suggests a significant disconnect between human-

perceived difficulty and model failure modes; problems that are not considered exceptionally hard for humans can still be robustly unsolvable by current models.

Benchmark Justification A strong benchmark should be realistic, difficult, and efficient (Nie et al., 2025). MATH-B meets these with (i) *realism*: problems drawn from standard high-school curricula; (ii) *difficulty*: base models perform poorly even with a large sampling budget (i.e. $\text{pass}@1024$); and (iii) *efficiency*: the small suite enables potentially rapid, low-cost evaluation, and—as shown in Section 4—post-trained models exhibit only marginal gains at $\text{pass}@1024$, making MATH-B a perfect test bed for research.

Benchmark Usage We expect researchers to use MATH-B with the base models and their corresponding splits specified in Section 3.2.3: (1) evaluate the base model to estimate its $\text{pass}@k$ for reasonably large k (e.g., $k = 1024$, expected to be nearly 0); (2) apply some RL method of interest; (3) re-estimate $\text{pass}@k$ for the post-trained policy, which indicates the *Expansion Rate* per Section 2. The benchmark is intended specifically for methods that aim to *expand* the listed base model’s reasoning boundary. While this is a research benchmark, we hope the methods that people come up with scales to larger scales and other scenarios

4 EVALUATING EXPANSION ACROSS FINETUNING METHODS

We evaluate several post-trained models on MATH-B to measure their ability to expand beyond their base model’s reasoning capabilities (i.e. Expansion Rate Section 2). Our analysis of models finetuned with reinforcement learning (RL) reveals that current methods achieve only modest expansion. As shown in Table 2, the three RL models based on DeepSeek-R1-Distill-Qwen1.5B (r1-1.5b) (DeepSeek-AI et al., 2025) solve fewer than 10% of the test problems. In contrast, Skywork-OR1-7B (skywork_or1) (He et al., 2025a) reaches a more promising 21% expansion; notably, its training involves adaptive entropy control and a higher temperature, likely affording a greater scope for exploration. Keep in mind that these are $\text{pass}@1024$ evaluations. This result suggests that RL techniques designed to explicitly encourage exploration may indeed result in higher Expansion Rates. We also observe that while prolonged training (i.e. more RL compute) can offer marginal gains—as seen in the slight improvement from Nemotron-Research-Reasoning-Qwen-1.5B version 1 (nemotron_v1) to Nemotron-Research-Reasoning-Qwen-1.5B version 2 (nemotron_v2) (Liu et al., 2025a)—the rather small 1.5% increase underscores the need for more efficient and effective exploration methods. See also the evolution of Expansion Rate in Figure 5.

Table 2: **Expansion Rates of post-trained models** using either RL or SFT/Distillation. The Expansion Rate measures the percentage of previously unsolvable problems (from the base model’s perspective) that the post-trained model can now solve We additionally add AIME24 ($\text{pass}@1$) numbers of the post-trained models to illustrate the difficulty of our dataset (He et al., 2025a; Yang et al., 2025; Liu et al., 2025a).

Base model	Post-trained	Method	Base unsolved	Expansion Rate (%, $\text{pass}@1024$)	AIME24 ($\text{pass}@1$)
<i>Reinforcement Learning (RL) models</i>					
r1-1.5b	nemotron_v1	RL	115	7.83	48.13
r1-1.5b	nemotron_v2	RL	115	9.57	49.58
r1-1.5b	DeepScaleR	RL	115	5.22	40.21
r1-7b	skywork_or1	RL	99	21.2	70.2
<i>Supervised Fine-Tuning (SFT) models</i>					
Qwen3-4B-base	Qwen3-4B	Long CoT Dist.	112	58.93	73.3
Qwen3-8B-base	Qwen3-8B	Long CoT Dist.	116	66.38	76.0

As an illustrative contrast, we evaluated Qwen3-4B and Qwen3-8B. These models are the result of finetuning their respective base versions (Qwen3-4B-Base and Qwen3-8B-Base) by distilling long Chain-of-Thought (CoT) reasoning trajectories from a more capable teacher model. They demonstrate substantially higher Expansion Rates of 58.93% and 66.38%, respectively. While this is not a direct apples-to-apples comparison due to differing training setups, the result is highly informative: it

shows that significant expansion is achievable when a model is exposed to the correct distribution of reasoning steps, an overlap that long CoT distillation provides (Yang et al., 2025).

This contrast highlights that the primary limitation of current RL techniques may not be an inherent inability of the base models to learn, but rather the failure of the exploration process to find these effective reasoning pathways on its own. Developing RL methods that can discover these pathways without a teacher model remains a key challenge for reaching frontier capabilities and is the central motivation for our work.

5 DISCUSSION AND RELATED WORK

On the Choice of $k = 1024$ Our selection of $k = 1024$ for all $\text{pass}@k$ evaluations is a deliberate choice to ensure our benchmark is challenging, stable, and efficient. Firstly, a large sampling budget is necessary to probe the true reasoning boundary of a model. Many popular benchmarks become saturated at this scale, with base models solving a high percentage of problems and leaving no room to measure improvement (Yue et al., 2025). We chose $k = 1024$ precisely because it represents a budget where MATH-B remains difficult, creating a meaningful testbed for genuine expansion. Secondly, this choice is strongly supported by our empirical analysis. While overall $\text{pass}@k$ performance shows a consistent log-linear increase with the sampling budget (Figure 6), the marginal gains for each additional sample diminish considerably (Figure 7). Most importantly, the Expansion Rate for RL-finetuned models—our core metric for progress—largely plateaus as the budget approaches 1024 (Figure 5). Therefore, $k = 1024$ represents a principled trade-off: it is large enough to push models beyond their comfort zone on a challenging benchmark, yet provides a stable and computationally feasible point to reliably measure the expansion of the reasoning boundary.

Limitations of Existing Benchmarks and Metrics In mathematical reasoning, progress is often measured by $\text{pass}@k$. For $k > 1$, this metric is taken to be indicative of a model’s exploratory potential. However, $\text{pass}@k$ is an incomplete measure of exploratory potential, as it conflates the sharpening of existing solutions (Consolidation) with the discovery of entirely new ones (Expansion) (Wu et al., 2025). As detailed in Section 2, our work is concerned with Expansion. Many existing benchmarks are now saturated by strong open-source base models, making it impossible to measure new boundary expansion (Yue et al., 2025; Balunović et al., 2025; Wu et al., 2025). Furthermore, these benchmarks are targets of hyper-optimization, potentially rewarding spurious correlations (Shao et al., 2025). MATH-B addresses this by serving as a diagnostic tool. It comprises problems that are topically standard high-school math but are constructed to expose the subtle failures and brittleness of the dominant open-source research paradigm.

Our Empirical Contribution in Context Building on the framework of Wu et al. (2025), we *instantiate* these ideas in a concrete, reusable benchmark. their work introduces and analyzes Expansion within a specific setting, whereas we put forth a zero-baseline dataset (i.e., $\text{pass}@1024 \approx 0$ for the base models) that operationalizes the concept and scales evaluation across a wide range of open-weight models. This offers a practical path to shift community focus from merely improving $\text{pass}@k$ on saturated suites to *demonstrably expanding the reasoning boundary*. Further, while MATH-B is instantiated for a specific set of open-weight base models, we expect a subset of items to remain zero-baseline for larger open-weight models. We did not evaluate those due to resource constraints. Since MATH-B is intended to drive RL methods that *expand* a given base model’s reasoning boundary (rather than chase model-specific quirks), we expect the resulting methods and insights to transfer across model families and scales.

6 CONCLUSION

Our findings indicate that current post-training methods in the open-source ecosystem, particularly for models up to 8B parameters, primarily refine pre-existing reasoning abilities rather than creating new ones. The poor performance of these post-trained models on MATH-B empirically confirms that they struggle to expand their capabilities to problems that lie just beyond their base model’s reach. We introduce MATH-B as a precise diagnostic tool to address this issue. Its purpose is to catalyze research into post-training methods that achieve genuine exploration, providing a clear and reliable signal for when a model has truly expanded the boundaries of machine reasoning.

7 REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our work. All code required to reproduce the experiments is provided in the supplementary material. Detailed derivations are included in the main text (Section 2). Experimental settings and hyperparameters are described in the main paper (Section 3.2) and supplementary sections (Appendix A.2).

REFERENCES

- Alon Albalak, Duy Phung, Nathan Lile, Rafael Rafailov, Kanishk Gandhi, Louis Castricato, Anikait Singh, Chase Blagden, Violet Xiang, Dakota Mahan, and Nick Haber. Big-math: A large-scale, high-quality math dataset for reinforcement learning in language models, 2025. URL <https://arxiv.org/abs/2502.17387>.
- Susan Amin, Maziar Gomrokchi, Harsh Satija, Herke Van Hoof, and Doina Precup. A survey of exploration methods in reinforcement learning. *arXiv preprint arXiv:2109.00157*, 2021.
- Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. Matharena: Evaluating llms on uncontaminated math competitions, February 2025. URL <https://matharena.ai/>.
- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- Zhipeng Chen, Xiaobo Qin, Youbin Wu, Yue Ling, Qinghao Ye, Wayne Xin Zhao, and Guang Shi. Pass@ k training for adaptively balancing exploration and exploitation of large reasoning models. *arXiv preprint arXiv:2508.10751*, 2025.
- Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei. Reasoning with exploration: An entropy perspective. *arXiv preprint arXiv:2506.14758*, 2025.
- Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. Go-explore: a new approach for hard-exploration problems, 2021. URL <https://arxiv.org/abs/1901.10995>.
- Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. Omni-math: A universal olympiad level mathematic benchmark for large language models, 2024a. URL <https://arxiv.org/abs/2410.07985>.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, et al. The language model evaluation harness, 07 2024b. URL <https://zenodo.org/records/12608602>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, et al. Olmo: Accelerating the science of language models, 2024. URL <https://arxiv.org/abs/2402.00838>.

- Nathan Habib, Cl  mentine Fourier, Hynek Kydl        , Thomas Wolf, and Lewis Tunstall. Lighteval: A lightweight framework for llm evaluation, 2023. URL <https://github.com/huggingface/lighteval>.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. Olympiad-bench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems, 2024. URL <https://arxiv.org/abs/2402.14008>.
- Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang, Jiacheng Xu, Wei Shen, Siyuan Li, Liang Zeng, Tianwen Wei, Cheng Cheng, Bo An, Yang Liu, and Yahui Zhou. Skywork open reasoner 1 technical report, 2025a. URL <https://arxiv.org/abs/2505.22312>.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning, 2025b. URL <https://arxiv.org/abs/2504.11456>.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL <https://arxiv.org/abs/2103.03874>.
- Hynek Kydl        . Math-verify: Math verification library, 2024. URL <https://github.com/huggingface/math-verify>.
- Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, September 2022. ISSN 1566-2535. doi: 10.1016/j.inffus.2022.03.003. URL <http://dx.doi.org/10.1016/j.inffus.2022.03.003>.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models, 2022. URL <https://arxiv.org/abs/2206.14858>.
- Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [<https://huggingface.co/AI-MO/NuminaMath-1.5>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.
- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models, 2025a. URL <https://arxiv.org/abs/2505.24864>.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025b.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing ol-preview with a 1.5b model by scaling rl. <https://tinyurl.com/KD2347>, 2025. Notion Blog.
- MAA. American invitational mathematics examination - aime. In *American Invitational Mathematics Examination - AIME*.
- MAA. American mathematics competitions. In *American Mathematics Competitions*, 2023.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013. URL <https://arxiv.org/abs/1312.5602>.

- Fan Nie, Ken Ziyu Liu, Zihao Wang, Rui Sun, Wei Liu, Weijia Shi, Huaxiu Yao, Linjun Zhang, Andrew Y. Ng, James Zou, Sanmi Koyejo, Yejin Choi, Percy Liang, and Niklas Muennighoff. Uq: Assessing language models on unsolved questions, 2025. URL <https://arxiv.org/abs/2508.17580>.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, et al. 2 olmo 2 furious, 2025. URL <https://arxiv.org/abs/2501.00656>.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, et al. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Negin Raoof, Etash Kumar Guha, Ryan Marten, Jean Mercat, Eric Frankel, Sedrick Keh, Hritik Bansal, Georgios Smyrnis, Marianna Nezhurina, Trung Vu, Zayne Rea Sprague, and Ludwig others. Evalchemy, 2025.
- Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei Du, Nathan Lambert, Sewon Min, Ranjay Krishna, et al. Spurious rewards: Rethinking training signals in rlvr. *arXiv preprint arXiv:2506.10947*, 2025.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–503, 2016. URL <http://www.nature.com/nature/journal/v529/n7587/full/nature16961.html>.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharmashan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm, 2017. URL <https://arxiv.org/abs/1712.01815>.
- Yuda Song, Julia Kempe, and Remi Munos. Outcome-based exploration for llm reasoning. *arXiv preprint arXiv:2509.06941*, 2025.
- Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning, 2024. URL <https://github.com/huggingface/trl>.
- Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language models with one training example. *arXiv preprint arXiv:2504.20571*, 2025.
- Fang Wu, Weihao Xuan, Ximing Lu, Zaid Harchaoui, and Yejin Choi. The invisible leash: Why rlvr may not escape its origin, 2025. URL <https://arxiv.org/abs/2507.14843>.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024. URL <https://arxiv.org/abs/2409.12122>.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, et al. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025a. URL <https://arxiv.org/abs/2503.14476>.

648 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian
649 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at
650 scale. *arXiv preprint arXiv:2503.14476*, 2025b.

651 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang.
652 Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?,
653 2025. URL <https://arxiv.org/abs/2504.13837>.
654

A APPENDIX

A.1 VERIFICATION ISSUES IN DIFFERENT FRAMEWORKS

Table 3: **Verification failure mode vulnerabilities across frameworks.** ✓ = vulnerable, − = resilient, ✓/− = partial. *Frameworks:* TRL (Transformers RL) (von Werra et al., 2024), VERL (Sheng et al., 2024), LM-Eval (hendrycks/minerva) (Gao et al., 2024b), LightEval (Habib et al., 2023), SCORE (LM-Eval SCORE math) (Gao et al., 2024b), evalchem (ZeroEval) (Raoof et al., 2025), HMMT (evalchem matharena) (Raoof et al., 2025), Math-V (Math-Verify) (Kydlicek, 2024). *Methods:* First# = first number in string; Right = rightmost priority; Ltd. = limited support (≤ 5 expressions). *F1–F7:* See Table 1 for detailed descriptions.

Failure Mode	TRL	VERL	LightEval	LM-Eval	SCORE	evalchem	HMMT	Math-V
F1: Multiple solutions (OR)	✓	✓	✓/−	✓	✓	✓	✓/−	✓/−
F2: Late correct	✓	✓/−	✓/−	✓/−	−	✓	−	✓/−
F3: Early correct	−	✓	✓	✓	−	−	✓	✓
F4: Answer corrections	✓	−	✓	−	−	✓	−	✓
F5: Unordered sets	✓	✓	✓	✓	✓	✓	−	✓
F6: Missing anchors	✓	−	−	−	−	−	✓	−
F7: MCQ partial match	✓	✓	✓	✓	✓	✓	✓	✓
Method	First	Last	Right	Last	All	First#	Last/All	Right
Anchors	Yes	No	No	No	No	No	No	Yes
Multi-Answer	No	No	Ltd.	No	Partial	No	Yes	Ltd.

A.2 INFERENCE PARAMETERS FOR EVALUATION

All models were run with nucleus sampling at a `top_p` of 0.95. Other parameters were set according to model-specific recommendations:

- **OLMo Models:** Temperature of 1.0, with `max_tokens` set to 2048 (v1) and 4096 (v2).
- **Qwen2.5-Math Models:** Temperature of 0.6 and `max_tokens` of 4096.
- **All Other Models:** Temperature of 0.6 and `max_tokens` of 32,768.

A.3 MATH-B DETAILS

In Figure 4, we show the distribution of final answers for MATH-B-U, which spans a broad range of integer values.

A.3.1 LOOKING AT SOME SAMPLES FROM MATH-B

In Tables 5 and 6, we list ten randomly sampled MATH-B-I questions—the base split and the full set, respectively.

A.4 ANALYSIS OF `PASS@K` PERFORMANCE SCALING

To further justify our evaluation methodology, we analyze the performance scaling of all 21 models from Table 4 on the MATH-B Union set.

In Figure 6, we plot the complete `pass@k` evolution as the sampling budget k increases up to 1024. The performance curves for nearly all models exhibit a characteristic log-linear growth, indicating that improvement is consistent but not linear with computational effort. While this trend suggests continued gains with more sampling, it also shows initial signs of plateauing at higher values of k .

To better quantify this observation, Figure 7 visualizes the marginal gain in performance. Specifically, it plots the average increase in the `pass@k` rate for each successive 64-sample increment. This plot clearly illustrates the principle of diminishing returns: the most significant gains are concentrated at lower sampling budgets, and the rate of improvement slows considerably as the budget approaches 1024. Together, these figures provide strong empirical support for our choice of $k = 1024$ as a

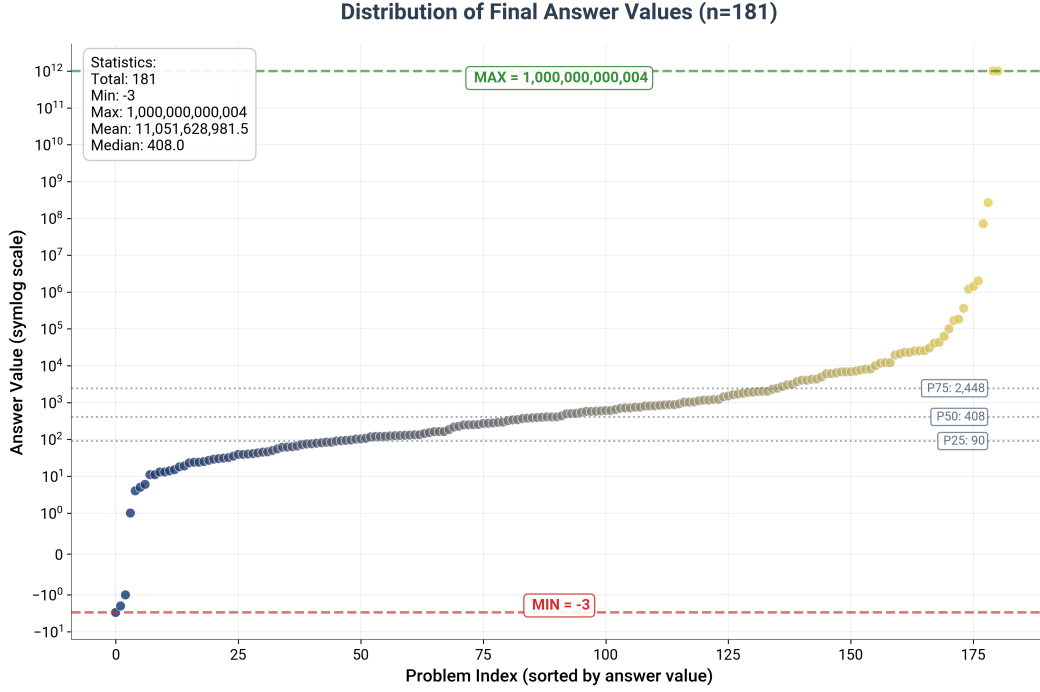


Figure 4: **Distribution of ground truth (final-answers) in MATH-B-U.** We use the log-scale for better readability.

practical and stable point for evaluation, beyond which brute-force sampling becomes an increasingly inefficient path to solving the remaining hard problems.

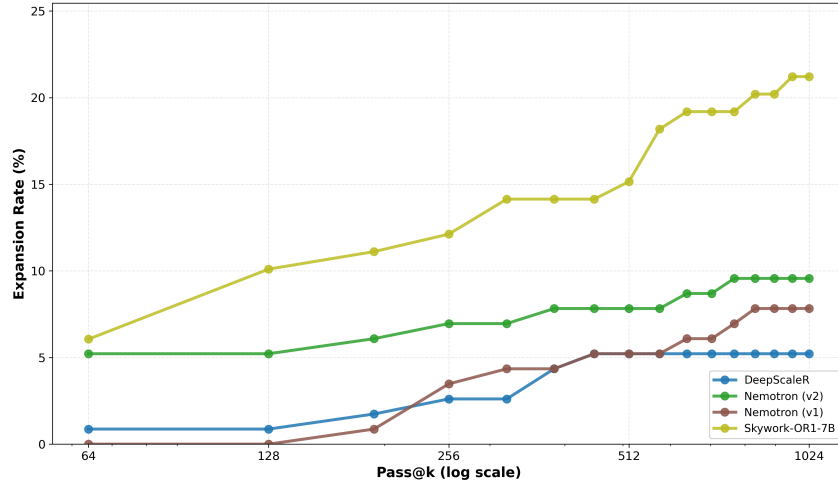


Figure 5: **Evolution of Expansion Rate for RL Models.** Models are evaluated on the MATH-B problems failed by their respective base models (115 for R1-Qwen2.5-1.5B; 99 for R1-Qwen2.5-7B).

A.5 LLM USAGE

The authors of this submission use IDEs with built-in LLM support, so LLMs have been used to help with menial coding tasks. Further, we used models like GPT-5, Gemini for rephrasing several

Table 4: **Models’ number of unsolved questions at pass@1024**. Models are grouped into Base vs. Supplementary. “Intersection (base)” and “Intersection (all)” indicate the overlap of unsolved problems across models.

Model	# Unsolved	Type	Notes
Qwen2.5-1.5B	145	Base	
Qwen2.5-7B	120	Base	
Qwen2.5-Math-1.5B	115	Base	
Qwen2.5-Math-7B	101	Base	
Qwen3-4B-Base	112	Base	
Qwen3-8B-Base	116	Base	
DeepSeek-R1-Qwen2.5-1.5B	115	Base	Distill
DeepSeek-R1-Qwen2.5-7B	99	Base	Distill
OLMo-7B	158	Base	
OLMo-2-7B	132	Base	
Llama-3.1-8B	151	Base	
Qwen2.5-1.5B-Instruct	124	Supplementary	
Qwen2.5-7B-Instruct	131	Supplementary	
Qwen2.5-Math-1.5B-Instruct	139	Supplementary	
Qwen2.5-Math-7B-Instruct	117	Supplementary	
Qwen3-4B	67	Supplementary	
Qwen3-8B	52	Supplementary	
DeepScaler-1.5B	142	Supplementary	
Nemotron-1.5B-v1	137	Supplementary	
Nemotron-1.5B-v2	142	Supplementary	
Skywork-OR1-7B	103	Supplementary	
Intersection (base)	41	—	Shared failures across all base models
Intersection (all)	13	—	Shared failures across all models

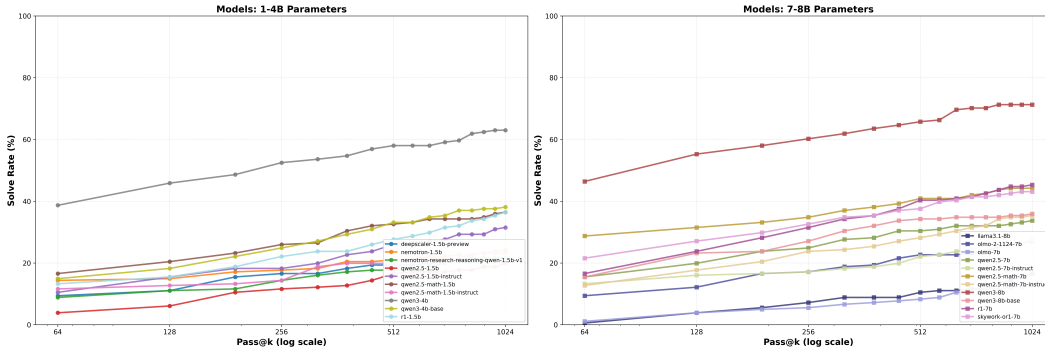


Figure 6: **Evolution of pass@k for all models on MATH-B-U.**

paragraphs in this manuscript. In addition to that, we used GPT-5-mini and o4-mini-high to label the difficulty and topics of the benchmark we create (Section 3.2).

Table 5: **Question texts (verbatim) from MATH-B-I (base)**. Each question includes the sentence “Let’s think step by step and output the final answer within `\boxed{ }`.” appended.

ID	Question (wrapped to fit)	Difficulty	Final Answer
1	Fisica and Ritmo discovered a piece of Notalium shaped like a rectangular box, and wanted to find its volume. To do so, Fisica measured its three dimensions using a ruler with infinite precision, multiplied the results and rounded the product to the nearest cubic centimeter, getting a result of 2017 cubic centimeters. Ritmo, on the other hand, measured each dimension to the nearest centimeter and multiplied the rounded measurements, getting a result of V cubic centimeters. Find the positive difference between the least and greatest possible positive values for V .	4.5	<code>7174</code>
2	Let \mathbb{N} denote the natural numbers. Compute the number of functions $f : \mathbb{N} \rightarrow \{0, 1, \dots, 16\}$ such that $f(x + 17) = f(x)$ and $f(x^2) \equiv f(x)^2 + 15 \pmod{17}$ for all integers $x \geq 1$.	6.0	<code>12066</code>
3	Each of the distinct letters in the following subtraction problem represents a different digit. Find the number represented by the word TEAM. PURPLE – COMET = TEAM.	3.5	<code>6852</code>
4	Find the sum of all positive integers n such that there exists an integer b with $ b \neq 4$ such that the base -4 representation of n is the same as the base b representation of n .	4.0	<code>1026</code>
5	Let S be a set of size 11. A random 12-tuple (s_1, \dots, s_{12}) of elements of S is chosen uniformly at random. Moreover, let $\pi : S \rightarrow S$ be a permutation of S chosen uniformly at random. The probability that $s_{i+1} \neq \pi(s_i)$ for all $1 \leq i \leq 12$ (where $s_{13} = s_1$) can be written as $\frac{a}{b}$ where a and b are relatively prime positive integers. Compute a .	6.0	<code>1000000000004</code>
6	Let \mathbb{N} denote the natural numbers. Compute the number of functions $f : \mathbb{N} \rightarrow \{0, 1, \dots, 16\}$ such that $f(x + 17) = f(x)$ and $f(x^2) \equiv f(x)^2 + 15 \pmod{17}$ for all integers $x \geq 1$.	4.5	<code>12066</code>
7	A unicorn is tethered by a 20-foot rope to the base of a cylindrical tower of radius 8 feet. The rope is attached to the tower at ground level and to the unicorn at a height of 4 feet. The rope is taut, its end is 4 feet from the nearest point on the tower, and the length of rope touching the tower is $\frac{a-\sqrt{b}}{c}$ feet, where a, b, c are positive integers and c is prime. Find $a + b + c$.	3.5	<code>813</code>
8	For which maximal N does there exist an N -digit number such that among any sequence of consecutive decimal digits some digit is present only once?	6.0	<code>1023</code>
9	Let $P = \{(x, y) \mid 0 \leq x, y \leq 25, x, y \in \mathbb{Z}\}$. Let T be the set of triangles formed by picking three distinct points in P (rotations, reflections, and translations count as distinct). Compute the number of triangles in T with area larger than 300.	1.5	<code>436</code>
10	Katie has a chocolate bar that is a 5-by-5 grid of square pieces, but she only wants to eat the center piece. She repeatedly (i) chooses a gridline and splits the bar, (ii) discards the part not containing the center, (iii) repeats until only the center piece remains. Compute the number of possible sequences of operations.	3.0	<code>6384</code>

Table 6: **Question texts (verbatim) from MATH-B-I (all)** set of all considered models in Section 3.2.3. Each question includes the sentence “Let’s think step by step and output the final answer within `\boxed{\hspace{1cm}}`” appended.

ID	Question (wrapped to fit)	Difficulty	Final Answer
1	Bob is writing a sequence of letters of the alphabet, each of which can be either uppercase or lowercase, according to the following two rules: If he had just written an uppercase letter, he can either write the same letter in lowercase after it, or the next letter of the alphabet in uppercase. If he had just written a lowercase letter, he can either write the same letter in uppercase after it, or the preceding letter of the alphabet in lowercase. For instance, one such sequence is $aAaABCDdcbBC$. How many sequences of 32 letters can he write that start at (lowercase) a and end at (lowercase) z ?	1.5	<code>\boxed{376}</code>
2	Let $V = \{1, \dots, 8\}$. How many permutations $\sigma : V \rightarrow V$ are automorphisms of some tree? (A graph consists of some set of vertices and some edges between pairs of distinct vertices. It is connected if every two vertices in it are connected by some path of one or more edges. A tree G on V is a connected graph with vertex set V and exactly $ V - 1$ edges, and an automorphism of G is a permutation $\sigma : V \rightarrow V$ such that vertices $i, j \in V$ are connected by an edge iff $\sigma(i)$ and $\sigma(j)$ are.)	5.5	<code>\boxed{30212}</code>
3	Find the sum of all positive integers n such that there exists an integer b with $ b \neq 4$ such that the base -4 representation of n is the same as the base b representation of n .	4.0	<code>\boxed{1026}</code>
4	A unicorn is tethered by a 20-foot silver rope to the base of a magician’s cylindrical tower whose radius is 8 feet. The rope is attached to the tower at ground level and to the unicorn at a height of 4 feet. The rope is taut, the end of the rope is 4 feet from the nearest point on the tower, and the length of the rope touching the tower is $\frac{a-\sqrt{b}}{c}$ feet, where a, b, c are positive integers and c is prime. Find $a + b + c$.	3.5	<code>\boxed{813}</code>
5	A 16x16 square sheet of paper is folded once in half horizontally and once in half vertically to make an 8x8 square. This square is again folded in half twice to make a 4x4 square. This square is folded in half twice to make a 2x2 square. This square is folded in half twice to make a 1x1 square. Finally, a scissor is used to make cuts through both diagonals of all the layers of the 1x1 square. How many pieces of paper result?	6.0	<code>\boxed{544}</code>
6	Matt writes a permutation of $\{1, 2, 3, \dots, 10\}$ across his paper with the leftmost number equal to 1 and the rightmost equal to 10. Exactly one interior number (not including 1 or 10) is less than both its immediate left and right neighbors. How many such permutations are there?	6.0	<code>\boxed{1636}</code>
7	Let \mathbb{N} denote the natural numbers. Compute the number of functions $f : \mathbb{N} \rightarrow \{0, 1, \dots, 16\}$ such that $f(x+17) = f(x)$ and $f(x^2) \equiv f(x)^2 + 15 \pmod{17}$ for all integers $x \geq 1$.	4.5	<code>\boxed{12066}</code>
8	Sir Alex plays the following game on a row of 9 cells. Initially all cells are empty. In each move he either (1) chooses any number of the form 2^j (nonnegative j) and puts it into an empty cell, or (2) chooses two cells with the same number 2^j , replaces the number in one cell with 2^{j+1} and erases the number in the other cell. At the end, one cell contains 2^n and the others are empty. Determine the maximum number of moves possible in terms of n . Provide the value when $n = 10$.	5.0	<code>\boxed{2025}</code>
9	Alice and Bob play on a board of one row of 2022 consecutive squares. They take turns placing domino tiles that cover two adjacent squares; Alice goes first. A tile must not cover a square already covered. The game ends when no tile can be placed. Alice wants to maximize the number of uncovered squares at the end; Bob wants to minimize it. What is the greatest number of uncovered squares Alice can ensure, no matter how Bob plays?	4.0	<code>\boxed{290}</code>
10	(Caos) A cao [sic] has 6 legs, 3 on each side. A walking pattern is an ordered sequence of raising and lowering each of the legs exactly once (total 12 actions), starting and ending with all legs on the ground. The pattern is safe if at any point he has at least 3 legs on the ground and not all three legs are on the same side. Estimate N , the number of safe patterns.	4.0	<code>\boxed{1416528}</code>

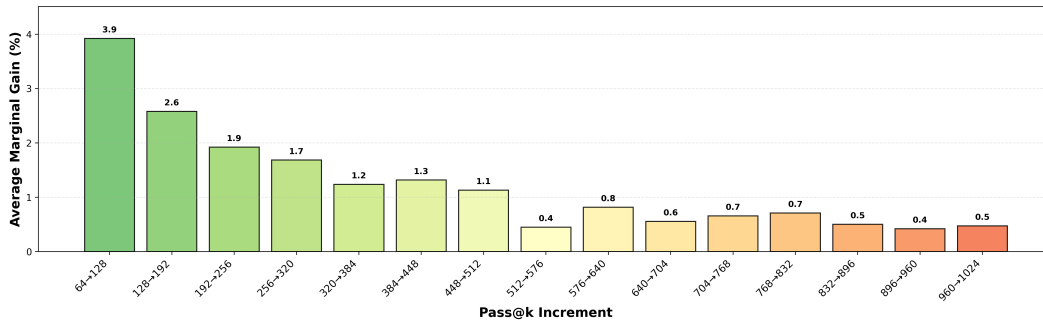


Figure 7: Average gains in **pass@k** relative to the size of MATH-B-U. Averaged over 21 models, the rate of solving new problems per 64-sample increment decreases as the total budget k grows, demonstrating diminishing returns.