

Scaling Up RL: Unlocking Diverse Reasoning in LLMs via Prolonged Training

NVIDIA

Recent advancements in reasoning-focused language models such as OpenAI’s O1 and DeepSeek-R1 have shown that scaling test-time computation-through chain-of-thought reasoning and iterative exploration—can yield substantial improvements on complex tasks like mathematics and code generation. These breakthroughs have been driven by large-scale reinforcement learning (RL), particularly when combined with verifiable reward signals that provide objective and grounded supervision. In this report, we investigate the effects of prolonged reinforcement learning on a small language model across a diverse set of reasoning domains. Our work identifies several key ingredients for effective training, including the use of verifiable reward tasks, enhancements to Group Relative Policy Optimization (GRPO), and practical techniques to improve training stability and generalization. We introduce controlled KL regularization, clipping ratio, and periodic reference policy resets as critical components for unlocking long-term performance gains. Our model achieves significant improvements over strong baselines, including +14.7% on math, +13.9% on coding, and +54.8% on logic puzzle tasks. To facilitate continued research, we release our model publicly.

Model on Hugging Face: [Nemotron-Research-Reasoning-Qwen-1.5B](#)

1. Introduction

The emergence of advanced reasoning models such as OpenAI’s O1 [1] and DeepSeek-R1 [2] represents a paradigm shift in Large Language Models (LLMs). This shift is characterized by scaling test-time computation-enabling models to perform extended deliberation through long-form Chain-of-Thought (CoT) reasoning. By allocating computational resources on using different reasoning strategies such as exploration, verification, and backtracking during the generation process, these models achieve remarkable improvements on complex tasks such as solving mathematical problems and code generation.

This paradigm shift has been primarily enabled by large-scale Reinforcement Learning (RL), which unlocks sophisticated reasoning capabilities beyond what supervised finetuning alone can achieve. The key innovation driving these advances is the use of tasks with verifiable rewards where correctness can be programmatically verified. This approach provides reliable training signals without the vulnerabilities inherent to learned reward models, which often suffer from reward hacking (such as generating superficially polite or unnecessarily verbose responses). Instead, verifiable rewards deliver grounded feedback that directly corresponds to correct reasoning and successful task completion [2, 3].

In this technical report, we investigate how prolonged reinforcement learning affects model performance, stability, and generalization across diverse reasoning tasks. Despite recent progress in this area, challenges remain in developing stable, scalable and sound training recipes for reasoning language models. These challenges include entropy collapse during training, difficulty in maintaining exploration-exploitation balance, and performance plateaus during extended optimization. Through systematic experimentation, we identify critical components that collectively form an effective framework for large-scale reinforcement learning of language models:

Diverse Training Data. We scale training across a wide variety of tasks with verifiable reward signals, spanning traditional reasoning domains such as mathematics and code generation, to more complex areas including STEM problem solving, logical puzzles, and instruction following. This data diversity helps evaluate the generality of reinforcement learning algorithms and exposes the model to a broader distribution of reasoning strategies. Compared to *DeepSeek-R1-Distill-Qwen-1.5B*, we achieve average pass@1 improvements of 14.7% on math benchmarks, 13.9% on coding, 54.8% on logic puzzles, 25.1% on STEM reasoning, and 18.1% on instruction-following.

Improvements over GRPO. Through extensive ablation studies, we incorporate and validate several

Data Type	Reward Type	Quantity	Data Source
Math	Binary	40k	DeepScaleR Dataset
Code	Continuous	24k	Eurus-2-RL Dataset
STEM	Binary	25k	SCP-116K Dataset
Logical Puzzles	Continuous	37k	Reasoning Gym
Instruction Following	Continuous	10k	Llama-Nemotron

Table 1 | Overview of training data used in our experiments, categorized by domain, reward type (binary or continuous), dataset size, and source. The datasets span a range of reasoning, coding, STEM, and instruction-following tasks, sourced from public repositories.

enhancements to Group Relative Policy Optimization (GRPO) proposed in DAPO [4], including decoupled clipping and dynamic sampling. These modifications lead to more efficient policy learning, with observed improvements as much as +5% on AIME2024 validation scores in our small-scale ablation study. In particular, dynamic sampling increases effective sample efficiency, and relaxed clipping coefficients over the importance sampling ratio contributes to the stabilization of entropy.

Training Stability. We found the training procedure to be sensitive to several critical hyperparameters, especially those related to sampling temperature, clipping thresholds, and KL divergence. While recent work suggests removing KL regularization, we find it beneficial in balancing exploration and exploitation, particularly when training from strong pretrained models. Applying a small but non-zero KL penalty, combined with entropy-preserving strategies such as those from DAPO, stabilizes learning and prevents entropy collapse. These adjustments enable scalable and stable training over hundreds of thousands of steps.

Prolonged Training. To support continued improvement beyond early convergence, we periodically reset the reference policy and optimizer states during training. Maintaining a static reference can restrict learning by penalizing beneficial divergence from the base model. Our results show that strategically resetting the reference policy, especially after observing KL spikes or validation performance declines, restores learning dynamics and enables further performance gains. For example, a mid-training reset recovered coding benchmark performance after a sharp decline, illustrating the effectiveness of this strategy in mitigating stagnation.

To support further research and development in this area, we are open-sourcing our model and hope that our contributions will be valuable to the broader research community.

2. Diverse Training Data

We scale training across a wide spectrum of tasks that provide verifiable reward signals as shown in Table 1. These tasks span from traditional reasoning domains, such as mathematical problem solving and code generation, to more complex and open-ended domains, including STEM-related problem solving, logical puzzles, and instruction following. The inclusion of such a diverse task set serves two key purposes. First, it broadens the model’s exposure to a wide distribution of reasoning patterns, encouraging generalization beyond narrow, domain-specific behaviors. This is especially critical for developing models of adapting to new or unseen task formulations. Second, the task diversity enables a more rigorous evaluation of RL algorithms, as it tests their ability to learn robust decision-making strategies across fundamentally different environments and reward structures. By grounding training in tasks with clear correctness criteria, we ensure that observed improvements are attributable to genuine progress in reasoning ability, rather than spurious correlations or overfitting to specific task formats.

2.1. Math

We use high-quality, community-curated datasets made available through DeepScaleR [3]. The training set consists of 40K math problems sourced from a diverse range of national and international math competitions. We adopt DeepScaleR’s original verifier and further augment it with an improved `math-verify`¹. We obtain the LLM’s answers by prompting the model with Let’s think step by step and output the final answer

¹<https://github.com/huggingface/Math-Verify>

within \boxed{\boxed{}}. We use a binary reward signal, assigning a score of 1 if the LLM’s response passes either the original or the enhanced `math-verify`, and 0 otherwise (for incorrect or improperly formatted answers).

2.2. Code

We utilize publicly available reinforcement learning datasets comprising 24K coding problems [5], sourced from various programming competitions. To support continuous reward feedback, we improve code execution environment to run all test cases rather than terminating on the first error and assign rewards based on the fraction of test cases passed. Submissions that fail to compile, contain syntax errors, or exceed a 5 second total timeout are assigned a reward of zero. We also include instructions for the LLM to enclose its final code response with triple backticks.

2.3. STEM

We use SCP-116K [6], a large-scale dataset containing 274K scientific problem-solution pairs spanning diverse fields such as physics, chemistry, biology, and mathematics. Each problem is accompanied by a corresponding solution extracted from the original source text, along with model-generated responses and reasoning paths produced by DeepSeek-R1. Given that SCP-116K was automatically extracted from heterogeneous and potentially noisy sources, we applied rigorous data filtering. First, we removed problems lacking a retrievable ground-truth solution from the source text. Then, we employed GPT-4o as a judge to assess whether the DeepSeek-R1 response aligned with the ground-truth answer. Only problems with consistent answers were retained, reducing the dataset from the original entries to 25K.

2.4. Logical Puzzles

The logical puzzles are well-suited for reasoning model training due to their broad coverage of different reasoning skills, as well as their clear objectives and evaluation metrics. We utilize the Reasoning Gym project [7], which offers approximately 100 tasks across various domains, including algebra, arithmetic, computation, cognition, geometry, graph theory, logic, and popular games. To facilitate model training and evaluation, we generate a large dataset consisting of 37K synthetic training samples and 9600 validation samples, spanning 96 tasks. Notably, some tasks have a unique solution, whereas others, such as the Rubik’s Cube and Countdown, admit multiple correct solutions. We employ the verifier provided by the Reasoning Gym repository for both model evaluation and reinforcement learning training signals. We use recommended default prompts which instruct models to enclose answers between `<answer>` `</answer>` tags.

2.5. Instruction Following

To enhance our model’s instruction-following capabilities, we leverage synthetic generated data [8] which data format is similar to IFEval [9]. Specifically the dataset contains synthetic prompts that pair tasks with randomly chosen instructions. For instance, a prompt may ask the model to “Write an essay about machine learning”, while the instruction specifies, “Your response should have three paragraphs.” We do not add further instructions on formatting and obtain the models response after thinking (`</think>` token). For evaluation, we adopt the strict mode of the IFEval verifier, which assesses only the original response and is considered more rigorous than the loose mode.

2.6. Implementation Details

To accommodate the complexity and diversity of our reward-generating tasks, we adopt a sandboxed reward server architecture. This design is motivated by several key factors. First, the diversity of data sources and task formats necessitates customized execution environments that can be tailored per task without interfering with the core training pipeline. Second, sandboxing allows us to isolate code execution from the training process, which is critical for security and fault-tolerance, especially when evaluating code generation or interacting with external systems. This ensures that any potential side effects or crashes during reward computation do not compromise the training run. Finally, we leverage multiprocessing and distributed reward servers across training clusters to scale reward evaluation efficiently. This parallelization is essential for maintaining high throughput during reinforcement learning, especially when reward signals to compute or

involve latency-bound operations such as code execution. We also launch reward calculations as soon as responses are completed, overlapping reward calculation with GPU tasks which do not depend on reward availability, such as $\pi_\theta(\tau), \pi_{old}(\tau)$ calculations.

3. Approach

We begin with a brief overview of the algorithm we used to optimize our language model, GRPO [10]. Our approach integrates techniques from DAPO [4], such as decoupled clipping and dynamic prompt sampling. Finally, to further address entropy collapse and training instability, we apply a KL divergence penalty and introduce a periodic reset of the reference policy. These components together enable stable training over extended durations and repeated epochs on the same dataset without degrading performance.

3.1. Background: Group Relative Policy Optimization

We adopt Group Relative Policy Optimization (GRPO) [10] as the core RL algorithm. Compared with Proximal Policy Optimization (PPO) [11], it removes the value model and instead use baseline estimates based on group scores. Formally the GRPO maximizes the following objective:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\min \left(r_\theta(\tau) A(\tau), \text{clip}(r_\theta(\tau), 1 - \epsilon, 1 + \epsilon) A(\tau) \right) \right], \quad (1)$$

where τ is the response sampled from the current policy π_θ . $r_\theta(\tau) = \frac{\pi_\theta(\tau)}{\pi_{old}(\tau)}$ is the probability ratio between the current policy and old policy before each actor update. The advantage used in GRPO foregoes the critic model of PPO, and instead estimates baseline from group scores $\{R_i\}_{i \in G(\tau)}$:

$$A(\tau) = \frac{R_\tau - \text{mean}(\{R_i\}_{i \in G(\tau)})}{\text{std}(\{R_i\}_{i \in G(\tau)})}. \quad (2)$$

3.2. Prolonged Reinforcement Learning

3.2.1. Mitigating Entropy Collapse

A key challenge in prolonged policy optimization is entropy collapse, a phenomenon where the model’s output distribution becomes overly peaked early in training, resulting in sharply reduced entropy. When entropy collapses, the policy prematurely commits to a narrow set of outputs, severely limiting exploration. This is particularly detrimental in methods like GRPO, where the learning signal depends on having a diverse set of sampled outputs to effectively estimate relative advantages. Without sufficient exploration, policy updates become biased, leading to stagnation in training.

A common mitigation strategy is to increase the sampling temperature during rollouts. However, we find that this approach merely delays the onset of entropy collapse rather than preventing it entirely, as entropy continues to decline steadily throughout training. As detailed in our temperature ablation study in Section 5.1, we used a high rollout temperature to promote exploration by boosting the initial entropy.

3.3. Decoupled Clip and Dynamic Sampling Policy Optimization (DAPO)

To address entropy collapse, we adopt several components from the DAPO algorithm [4], which are specifically designed to maintain exploration and output diversity. First, DAPO introduces decoupled clipping, where the lower and upper clipping bounds in the PPO objective are treated as separate hyperparameters:

$$\text{clip}(r_\theta(\tau), 1 - \epsilon_{low}, 1 + \epsilon_{high}). \quad (3)$$

By setting a higher value for ϵ_{high} , the algorithm promotes ‘clip-higher’, uplifting the probabilities of previously unlikely tokens and encouraging broader exploration. We find that this modification helps retain entropy and reduces premature mode collapse.

Additionally, DAPO employs dynamic sampling, filtering out prompts for which the model consistently succeeds or fails (i.e., accuracy 1 or 0), as these provide no learning signal. This focus on intermediate difficulty examples further helps maintain a diverse and stable learning signal during training.

3.3.1. KL Regularization and Reference Policy Reset

While DAPO and temperature adjustment help slow entropy collapse, we find that explicit regularization via a KL divergence penalty provides a stronger and more stable solution. Specifically, we incorporate a KL penalty between the current policy π_θ and a reference policy π_{ref} :

$$L_{KL-RL}(\theta) = L_{GRPO}(\theta) - \beta D_{KL}(\pi_\theta || \pi_{ref}), \quad (4)$$

where the following unbiased estimator is commonly used[12]:

$$D_{KL}(\pi_\theta || \pi_{ref}) = \frac{\pi_{ref}(\tau)}{\pi_\theta(\tau)} - \log \frac{\pi_{ref}(\tau)}{\pi_\theta(\tau)} - 1. \quad (5)$$

This penalty not only helps maintain entropy but also serves as a regularizer to prevent the online policy from drifting too far from a stable reference, stabilizing learning and mitigating overfitting to spurious reward signals.

Recent works [4, 13, 14, 15] have argued for the removal of the KL penalty, citing that models naturally diverge during training on chain-of-thought reasoning tasks. We observe that this perspective often applies to base models prior to any supervised fine-tuning. In contrast, we begin from a well-initialized checkpoint (DeepSeek-R1-Distill-Qwen-1.5B) already capable of generating coherent CoT outputs. In this context, retaining a KL penalty is still beneficial for both stability and sustained entropy.

We further observe that as training progresses, the KL term may increasingly dominate the loss, leading to diminishing policy updates. To alleviate this, we introduce a simple yet effective technique: *reference policy reset*. Periodically, we hard-reset the reference policy π_{ref} to a more recent snapshot of the online policy π_θ , and reinitialize the optimizer states. This allows the model to continue improving while maintaining the benefits of KL regularization. We apply this reset strategy throughout training to avoid premature convergence and encourage prolonged training.

4. Experiment Results

To ensure stable and effective reinforcement learning across diverse reasoning tasks, we adopt a staged training strategy involving multiple sequential runs, each designed to address specific challenges observed in earlier phases. This approach allows us to iteratively refine model behavior, incorporate additional data sources, adjust hyperparameters, and reset training dynamics when necessary. Throughout, we closely monitor key training signals, including KL divergence, entropy, and response length, and scores on a held-out validation set, using them as indicators to guide interventions such as reward shaping or context window changes. Our complete training recipe reflects this progressive, intervention-driven process, which we detail below.

4.1. Training Setup

We primarily leverage open-source framework verl [16] for reinforcement learning training. We adopt proposed enhancements from DAPO [4], decoupling clipping hyperparameters with $\epsilon_{low} = 0.2$, $\epsilon_{high} = 0.4$, and dynamic sampling for filtering prompts with accuracy equal to 1 and 0. We apply KL divergence penalty with a small penalty coefficient $\beta = 0.0001$. For rollout, we sample $n = 16$ responses for each prompt with a context window limit of 8096 and use a high sampling temperature of 1.2. We set batch size to 256 and mini-batch size to 64 (equating to 4 gradient updates per rollout step). For training we use the AdamW [17] optimizer with a constant learning rate of 2×10^{-6} . We use *DeepSeek-R1-Distill-Qwen-1.5B* as the initial base model policy. We conduct training on 4 8 x NVIDIA-H100-80GB nodes, and the whole training runs for approximately 16k GPUs hours.

4.2. Training Recipe

We construct a validation data blend to closely monitor training progress across steps. This validation set includes subsets from our evaluation benchmark, specifically AIME2024, Codeforces, GPQA-diamond, IFEval, and the logic puzzle *graph_color* from Reasoning Gym. We evaluate model performance using similar sampling parameters as in evaluation settings other than the same context window used in training. Occasionally, we

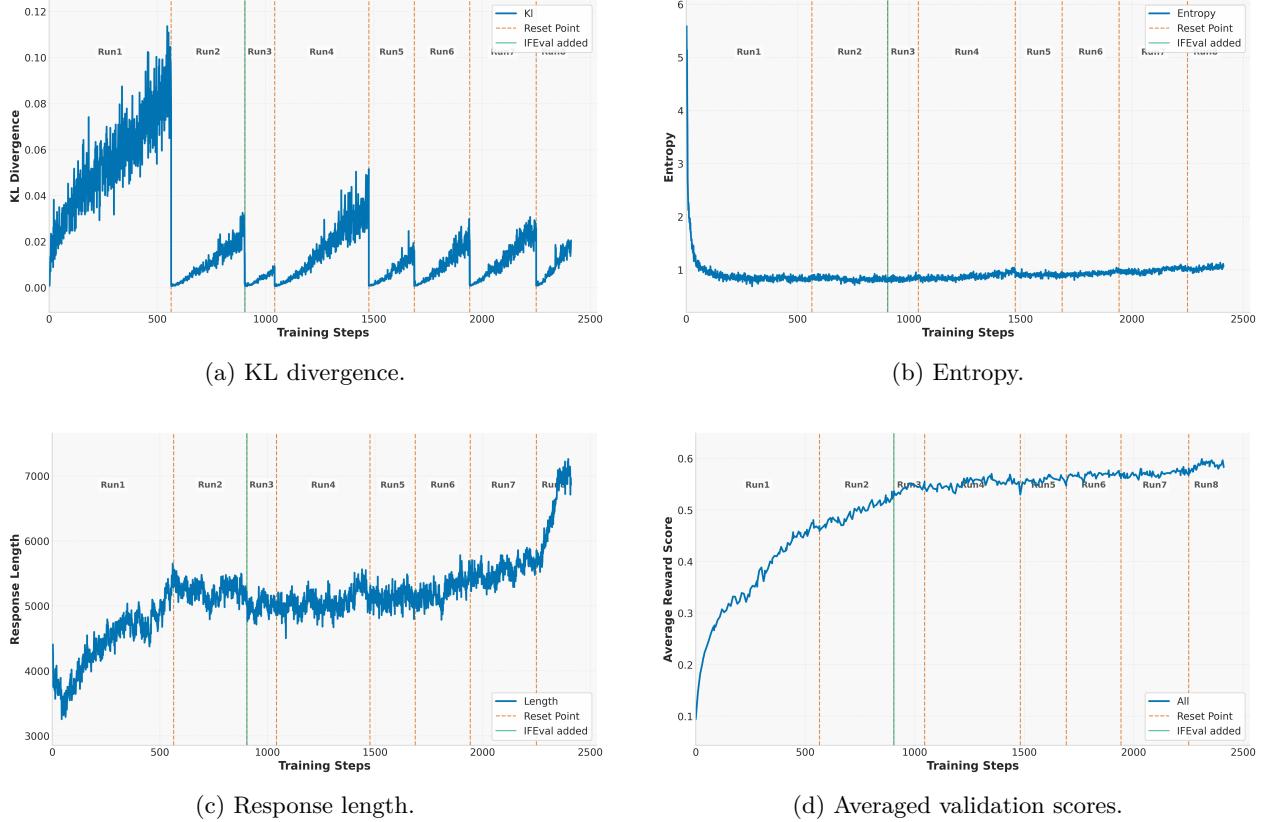


Figure 1 | Training dynamics across multiple runs. This figure presents key metrics monitored during training, including KL divergence, entropy, response length, and validation performance across different stages. Each run reflects specific adjustments to training conditions, such as changes in context window size, reward shaping, and rollout counts. The dynamics highlight the impact of these adjustments on model stability, response quality, and performance improvements on validation scores.

perform a hard reset of the reference model, as described in Section 3.3.1, particularly when validation metrics significantly degrade or when improvements plateau. Unless otherwise specified, we use the hyperparameters detailed in Section 4.1.

Interestingly, the hard reset (resetting reference policy and optimizer states) not only restores training stability but also provides an opportunity to adjust training hyperparameters and introduce enhancements such as additional training data and reward shaping. Figure 1 presents key statistics on training dynamics across different stages. The final training recipe comprises several sequential stages, described below:

- **Run 1:** We begin training on four tasks from Section 2, excluding instruction-following data. In this phase, the model initially adapts to a reduced context window of 8k, leading to shorter responses. Subsequently, response length increases along with improved validation scores. Toward the end of this stage, we observe instability and degradation in validation performance.
- **Run 2:** We perform a hard reset of the reference policy and resume training with the same setup as Run 1. Unlike DeepScaleR [3], which proposes increasing the context window, we maintain a fixed 8k context. Notably, validation scores continue to improve in this stage, even as response length remains relatively stable.
- **Run 3:** We incorporate instruction-following data into the training mix and continue training. This stage proceeds until we observe a sudden increase in response length, primarily due to the model repeating answers and failing to terminate with an *eos* token.
- **Runs 4 and 5:** We introduce reward shaping by penalizing responses that do not terminate correctly.

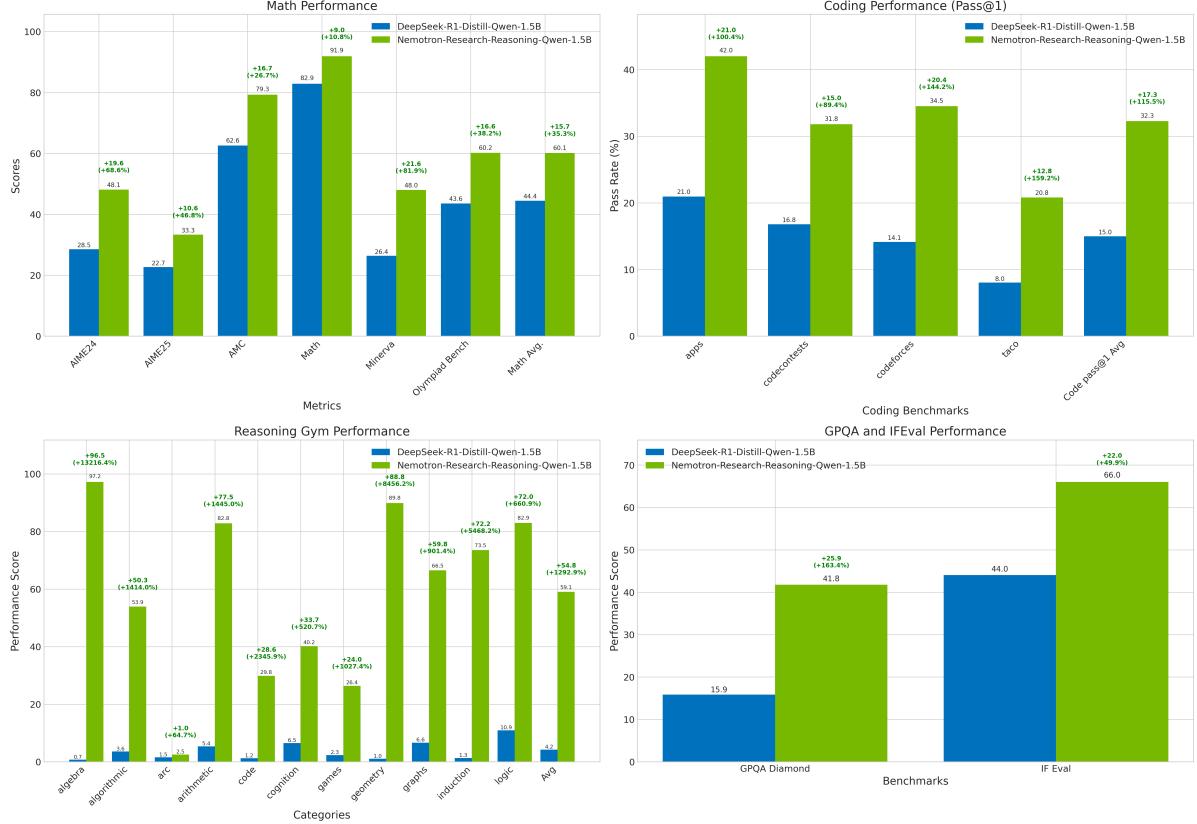


Figure 2 | Performance comparison between *DeepSeek-R1-Distill-Qwen-1.5B* and our model *Nemotron-Research-Reasoning-Qwen-1.5B* across multiple benchmarks in diverse domains, including math, coding, Reasoning Gym puzzles, GPQA Diamond, and IFEval.

This encourages proper generation behavior, resulting in a modest reduction in response length. Toward the end of this stage, gains on validation benchmarks begin to plateau.

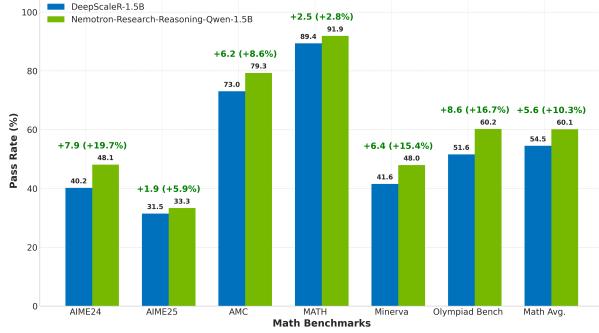
- **Runs 6 and 7:** We increase the rollout count from 16 to 32, performing two hard resets in the process. Interestingly, response length begins to rise again alongside improvements in validation metrics.
- **Run 8:** We extend the context window to 16k tokens and reduce rollout count to 16. Despite the model being trained on an 8k context window for most of the time, it quickly adapts to the extended context window. We observe marginal improvements in hard math tasks like AIME, with more substantial gains coming from other domains.

4.3. Evaluation Benchmarks

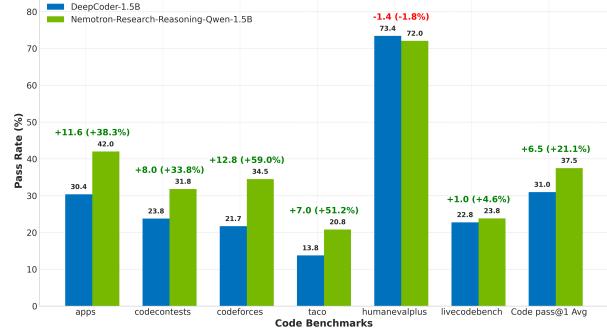
We evaluate models on the breadth of various tasks across math, coding, reasoning, and instruction following. For math, we follow DeepScaleR [3] and SimpleRL [18], and evaluate on AIME2024 [19], AIME2025 [20], AMC [21] (composed of AMC2022 and AMC2023), MATH [22], Minerva Math [23], and Olympiad Bench [24]. For coding, we use the validation set from PRIME [5] consisted of APPS [25], Codecontests [26], Codeforces², and TACO [27]. We also include benchmarks HumanevalPlus [28] and LiveCodeBench [29]. For logic puzzles, we reserved 100 samples from each reasoning gym [7] tasks as test datasets for evaluation. In addition, we use a curated subset³ from GPQA Diamond [30] and IFEval [31] to evaluate the capability of our models in STEM reasoning and instruction following [9].

²<https://huggingface.co/datasets/MatrixStudio/Codeforces-Python-Submissions>

³<https://huggingface.co/datasets/spawn99/GPQA-diamond-ClaudeR1>



(a) Comparision with DeepScaleR-1.5B [3].



(b) Comparision with DeepCoder-1.5B [13].

Figure 3 | Performance comparison of *Nemotron-Research-Reasoning-Qwen-1.5B* with domain-specialized models *DeepScaleR-1.5B* and *DeepCoder-1.5B*. Trained on diverse-domain data, our model achieves competitive results in both mathematical reasoning and code generation tasks.

4.4. Evaluation Results

We use vllm (version 0.7.2) [32] for inference backend, with sampling temperature of 0.6, nucleus sampling [33] with $\text{top_p} = 0.95$ and a context windo of 32k. For math, coding and STEM reasoning tasks we obtain estimates of *pass@1* from 16 samples for each benchmark prompts from strictly binary rewards. For other tasks such as logical puzzles and instruction following, we calculate the average continuous reward score from our rule-based verifiers with 16 samples. We evaluate and report benchmark results for open-source models using our own evaluation settings. We acknowledge that discrepancies from originally reported numbers may arise due to differences in evaluation settings and inherent variability from inference.

Figure 2 provides a detailed comparison between *DeepSeek-R1-Distill-Qwen-1.5B* and our model, *Nemotron-Research-Reasoning-Qwen-1.5B*, across multiple domains: math, coding, reasoning, and instruction following. In the math domain, our model consistently outperforms across benchmarks, showing an average improvement of 15.7%, which reflects enhanced symbolic and numerical reasoning capabilities. In coding, *Nemotron-Research-Reasoning-Qwen-1.5B* surpasses *DeepSeek-R1-Distill-Qwen-1.5B* in competitive programming tasks as measured by *pass@1* accuracy, indicating superior code synthesis ability. Our model also demonstrates substantial gains in STEM reasoning and instruction following, with improvements of 25.9% on GPQA Diamond and 22.0% on IFEval. For logic puzzles in the Reasoning Gym suite, we adopt the categorization of 96 tasks as defined by the official GitHub repository. Notably, *DeepSeek-R1-Distill-Qwen-1.5B* underperforms even on relatively simple mathematical tasks such as algebra and arithmetic. Closer inspection reveals that the model consistently formats its answers using `\boxed{}` rather than adhering to the instruction to use `<answer></answer>` tags. Despite poor initial formatting behavior, the model is able to achieve high accuracy on these easier tasks post training, suggesting that formatting is relatively easy to learn. Our models still exhibit room for improvement on more challenging categories, including tasks from arc, code, cognition, and games. In these cases, the model often fails to make meaningful progress. Further analysis indicates that these failures stem from either a lack of core reasoning skills necessary to solve specific subtasks or insufficient background knowledge related to the problem domains. Addressing these limitations may require additional finetuning data to better support model from a cold start, which we leave these enhancements to future work.

We further compare the performance of *Nemotron-Research-Reasoning-Qwen-1.5B* against two domain-specialized baselines, *DeepScaleR-1.5B* (optimized for mathematical reasoning) and *DeepCoder-1.5B* (focused on competitive programming tasks) in Figure 3. Trained on a broad and diverse set of domains rather than being narrowly specialized, our model demonstrates strong generalization capabilities and achieves competitive performance across both math and code benchmarks.

5. Ablation Studies

To better understand the contribution of key components in our training pipeline, we conduct a series of ablation studies. Specifically, we examine the effects of rollout temperature sampling, decoupled clipping and

dynamic sampling enhancements proposed in DAPO [4]. Additionally, we analyze the training stability, i.e., performance degrades or plateaus over extended runs. These controlled variations help isolate the impact of each mechanism and offer insights into both the robustness and limitations of our current setup.

5.1. Rollout Sampling Temperature



Figure 4 | Ablation study on rollout sampling temperature during early- (Run 1) and late-stage (Run 7) training.

We conduct ablations on the rollout sampling temperature to study its effect on training dynamics and keep performance metrics shown in Figure 4. Sampling temperature plays a crucial role in shaping the distribution of rollout candidates, thereby influencing both reward signal quality and model learning behavior. Our study includes two separate evaluations: one during early-stage training (Run 1) and another during late-stage training (Run 7), before extending the context window.

Early-stage Training Ablation (Run 1) In the early reinforcement learning phase, we experiment with different sampling temperatures. We find that training with a low temperature (e.g., 0.6) leads to early instability. Both the average reward and downstream performance on AIME2024 degrade significantly, likely due to reduced diversity in rollouts and overexploitation of suboptimal generation patterns. In contrast, using a higher temperature of 1.2 yields more stable training. Although initial rewards are lower compared to temperature 0.6, they improve steadily and eventually surpass the low-temperature setting. The accuracy score of AIME2024 follows a similar trend to the reward as the average reward in term of sampling temperature, indicating that more diverse rollouts support better generalization and learning at this stage.

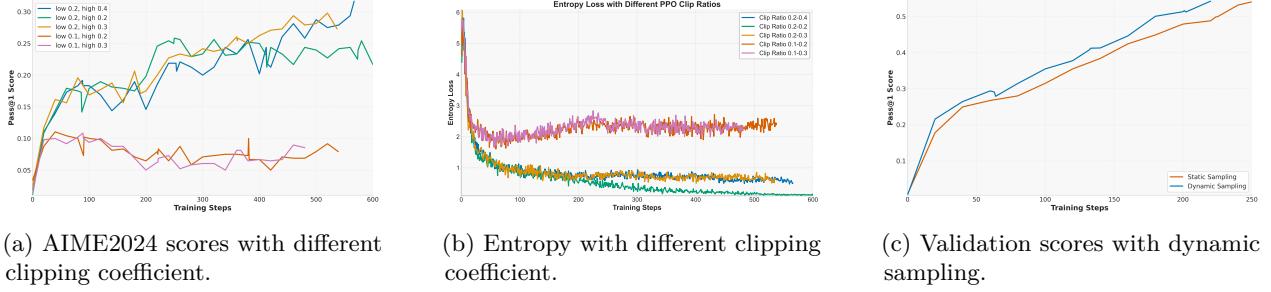
Late-stage Training Ablation (Run 7) We further evaluate temperature settings during Run 7. After extensive training and increased rollout counts, the model is already strong. Conventional intuition might suggest that lower temperatures would be more effective in this context to leverage more exploitation over exploration. However, our ablation shows that high temperatures again yield better reward progression and validation stability than lower temperatures. The high-temperature rollouts prevent mode collapse and support continued progress, even when performance gains begin to plateau.

These results underscore the importance of temperature tuning throughout training. Lower temperatures may prematurely limit exploration, especially early in training, while properly higher temperatures enable broader behavioral exploration and improved reward optimization. We note that this observation may be model-dependent, and transferring it to other models may require further analysis.

5.2. Decoupled Clip and Dynamic Sampling

We investigate the impact of decoupled clipping by varying its lower and upper thresholds, ϵ_{low} and ϵ_{high} , as illustrated in Figure 5. Our results show that setting $\epsilon_{low} = 0.2$ outperforms more conservative choices like 0.1, indicating that more aggressively down-weighting actions with negative advantage can facilitate learning. Additionally, increasing the upper threshold to $\epsilon_{high} = 0.4$ led to further improvements, mitigating entropy collapse and yielding the highest validation performance overall.

Interestingly, we observed that applying a smaller ϵ_{low} led to higher policy entropy. This is because a smaller ϵ_{low} discourages the policy from strongly penalizing actions with negative advantage, effectively



(a) AIME2024 scores with different clipping coefficient.

(b) Entropy with different clipping coefficient.

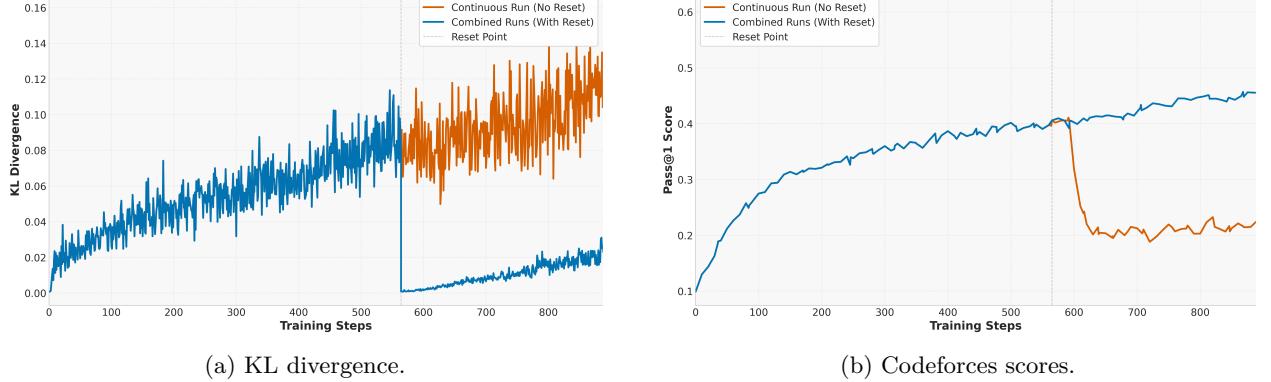
(c) Validation scores with dynamic sampling.

Figure 5 | Ablation study on decoupled clipping and dynamic sampling [4]. Both techniques demonstrate clear improvements over standard GRPO [10] in our experiments.

flattening the action distribution and encouraging broader exploration. However, this increased entropy did not translate into improved learning outcomes, emphasizing that excessive exploration alone is insufficient, and that a careful balance between exploration and exploitation is crucial for stable and effective training.

Finally, we verified that dynamic sampling led to faster improvements than static sampling by filtering out prompts with zero advantage. This increases reward signal density in each batch, thereby improving the overall sample efficiency of training.

5.3. Resetting Reference Policy



(a) KL divergence.

(b) Codeforces scores.

Figure 6 | Extended training sometimes becomes unstable. In our case, extending Run 1 causes a sudden degradation on Codeforces validation scores. In other cases, improvement plateaus with a sudden increase in KL divergence. We perform a hard reset of training, resetting the reference policy used for KL divergence calculation as well as optimizer states.

We observe that extending training with the same setting and training hyperparameters does not always lead to consistent improvements. As shown in Figure 6, if we simply extend the training of Run 1, it encounters a sharp drop in validation scores on the Codeforces benchmark. In other scenarios, performance gains plateau and the KL divergence against the reference policy spikes unexpectedly. To mitigate these issues, we implement a hard reset strategy, refreshing both the reference policy (used for KL divergence computation) and the optimizer states, to restore stability and enable further effective training.

5.4. Mitigating Entropy Collapse

Reinforcement learning often suffers from entropy collapse, where the model’s output distribution becomes too concentrated early in training, leading to a rapid drop in entropy. This can hinder exploration and bias advantage estimation, ultimately stalling learning. Figure 7 illustrates the phenomenon: without intervention, the model’s output entropy drops sharply early in training and remains low, signaling a premature convergence to narrow output distributions. To better understand and mitigate entropy collapse during reinforcement

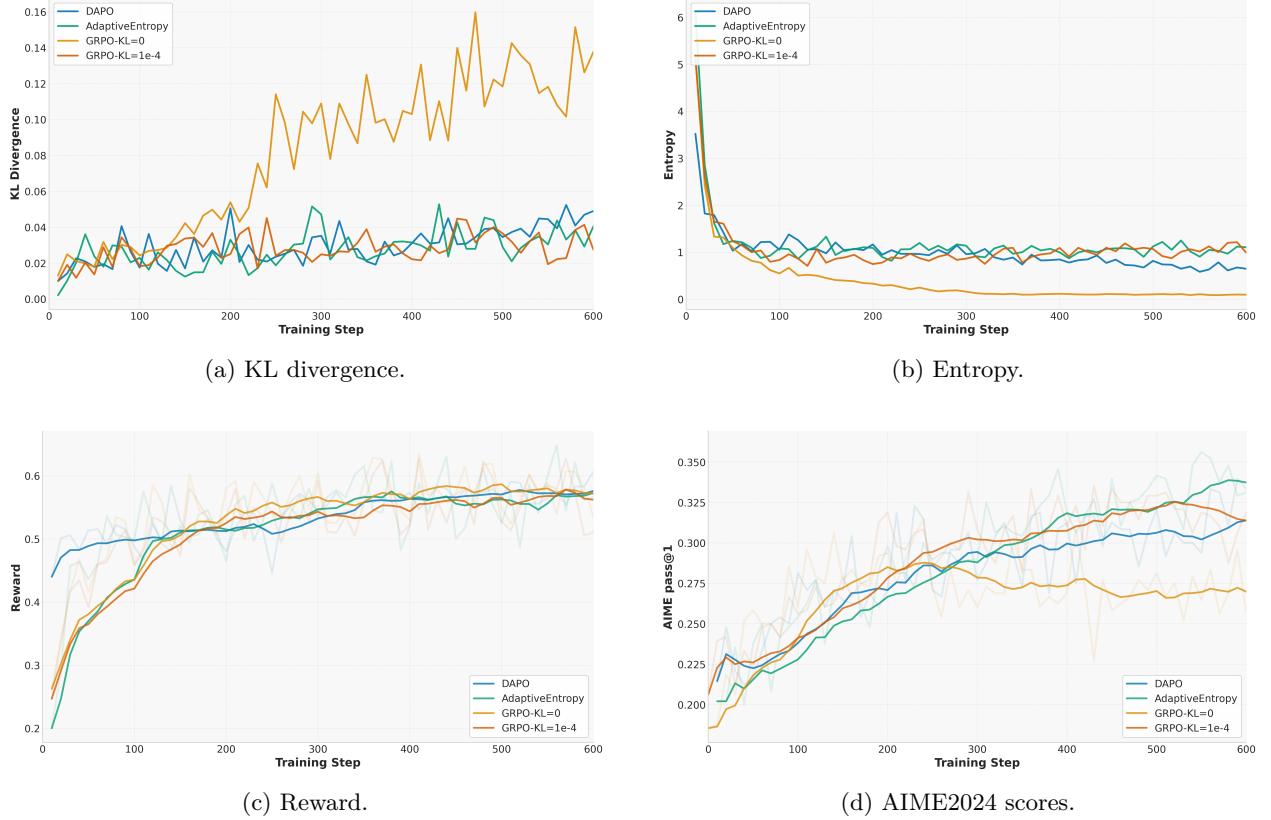


Figure 7 | Training dynamics under different entropy mitigation strategies.

learning, we conducted a series of ablation studies evaluating the training dynamics under different strategies. We investigated the following settings:

- **GRPO-KL=0.** We employ standard GRPO training without KL divergence penalty or entropy bonus.
- **GRPO-KL=1e-4** We apply a small KL regularization term to control the divergence between the online policy and a reference policy. This helps prevent the model from drifting too far and encourages continued exploration.
- **DAPO.** We adopt the ‘clip-higher’ technique from DAPO [4], which upweights the probabilities of lower-ranked tokens during sampling.
- **AdaptiveEntropy.** Inspired by parallel work [15], we also experiment with an adaptive entropy strategy, where the entropy bonus coefficient is dynamically enabled or disabled based on training statistics. Since tuning a fixed entropy coefficient is notoriously brittle, this approach sets a target entropy and selectively applies the bonus to maintain a desired level of output diversity.

We observe that all three methods are effective in mitigating entropy collapse to varying degrees. The adaptive entropy strategy introduces the most flexibility, but requires careful calibration of the target entropy and adds one more sensitive hyperparameter to tune. In contrast, DAPO and KL penalty offer a more conservative and robust solution, especially when training from a capable base model. Based on these findings, we adopt a combination of DAPO and a KL penalty in our final training setup. The KL term not only contributes to entropy preservation but also improves training stability by preventing runaway divergence from the reference model. Together, these methods strike a balance between preserving exploration and maintaining coherent, high-quality generations throughout prolonged training.

6. Conclusion

In this technical report, we presented a comprehensive investigation of prolonged reinforcement learning for reasoning-focused language models, identifying critical components that enable stable and effective training across diverse tasks. Our work demonstrates that through careful algorithm design, including decoupled clipping, dynamic sampling, controlled KL regularization, and periodic reference policy resets, even small-scale models can achieve substantial reasoning improvements without the computational demands of larger architectures. These techniques yielded significant performance gains over the baseline DeepSeek-R1-Distill-Qwen-1.5B model across multiple domains, including mathematics (+14.7%), coding (+13.9%), logic puzzles (+54.8%), STEM reasoning (+25.1%), and instruction-following (+18.1%). Importantly, our approach maintained competitive performance with domain-specialized models despite being trained on a broader range of tasks, suggesting that properly implemented reinforcement learning and prolonged training can effectively bridge the gap between general-purpose models and specialized reasoning systems. By open-sourcing our model and sharing our training methodology, we hope to facilitate further advancements in alignment, optimization, and reasoning within resource-efficient language models.

A. Contributors

Mingjie Liu, Shizhe Diao, Jian Hu, Ximing Lu, Xin Dong, Hao Zhang, Alexander Bukharin, Shaokun Zhang, Jiaqi Zeng, Makesh Narsimhan Sreedhar, Gerald Shen, David Mosallanezhad, Di Zhang, Jonas Yang, June Yang, Oleksii Kuchaiev, Guilin Liu, Zhiding Yu, Pavlo Molchanov, Yejin Choi, Jan Kautz, Yi Dong

References

- [1] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [2] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [3] 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 o1-preview with a 1.5b model by scaling rl. <https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>, 2025. Notion Blog.
- [4] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, 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, Weinan Dai, 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, 2025.
- [5] Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- [6] Dakuan Lu, Xiaoyu Tan, Rui Xu, Tianchu Yao, Chao Qu, Wei Chu, Yinghui Xu, and Yuan Qi. Scp-116k: A high-quality problem-solution dataset and a generalized pipeline for automated extraction in the higher education science domain, 2025.
- [7] Zafir Stojanovski, Oliver Stanley, Joe Sharratt, Richard Jones, Abdulhakeem Adefioye, Jean Kaddour, and Andreas Köpf. Reasoning gym: Reasoning environments for reinforcement learning with verifiable rewards, 2025.
- [8] Nvidia, :, Bo Adler, Niket Agarwal, Ashwath Aithal, Dong H. Anh, Pallab Bhattacharya, Annika Brundyn, Jared Casper, Bryan Catanzaro, Sharon Clay, Jonathan Cohen, Sirshak Das, Ayush Dattagupta, Olivier Delalleau, Leon Derczynski, Yi Dong, Daniel Egert, Ellie Evans, Aleksander Ficek, Denys Fridman, Shaona Ghosh, Boris Ginsburg, Igor Gitman, Tomasz Grzegorzek, Robert Hero, Jining Huang, Vibhu Jawa, Joseph Jennings, Aastha Jhunjhunwala, John Kamalu, Sadaf Khan, Oleksii Kuchaiev, Patrick LeGresley, Hui Li, Jiwei Liu, Zihan Liu, Eileen Long, Ameya Sunil Mahabaleshwarkar, Somshubra Majumdar, James Maki, Miguel Martinez, Maer Rodrigues de Melo, Ivan Moshkov, Deepak Narayanan, Sean Narendhiran, Jesus Navarro, Phong Nguyen, Osvald Nitski, Vahid Noroozi, Guruprasad Nutheti, Christopher Parisien, Jupinder Parmar, Mostofa Patwary, Krzysztof Pawelec, Wei Ping, Shrimai Prabhuloye, Rajarshi Roy, Trisha Saar, Vasanth Rao Naik Sabavat, Sanjeev Satheesh, Jane Polak Scowcroft, Jason Sewall, Pavel Shamis, Gerald Shen, Mohammad Shoeybi, Dave Sizer, Misha Smelyanskiy, Felipe Soares, Makesh Narsimhan Sreedhar, Dan Su, Sandeep Subramanian, Shengyang Sun, Shubham Toshniwal, Hao Wang, Zhilin Wang, Jiaxuan You, Jiaqi Zeng, Jimmy Zhang, Jing Zhang, Vivienne Zhang, Yian Zhang, and Chen Zhu. Nemotron-4 340b technical report, 2024.
- [9] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023.
- [10] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [11] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- [12] John Schulman. Approximating KL Divergence.
- [13] Michael Luo, Sijun Tan, Roy Huang, Xiaoxiang Shi, Rachel Xin, Colin Cai, Ameen Patel, Alpay Ariyak, Qingyang Wu, Ce Zhang, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepcoder: A fully open-source 14b coder at o3-mini level. <https://pretty-radio-b75.notion.site/DeepCoder-A-Fully-Open-Source-14B-Coder-at-03-mini-Level-1cf81902c14680b3bee5eb349a512a51>, 2025. Notion Blog.
- [14] Yu Yue, Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, Xiangpeng Wei, Xiangyu Yu, Gaohong Liu, Juncai Liu, Lingjun Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Ru Zhang, Xin Liu, Mingxuan Wang, Yonghui Wu, and Lin Yan. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks, 2025.

- [15] Jujie He, Jiacai 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 series. <https://capricious-hydrogen-41c.notion.site/Skywork-Open-Reasoner-Series-1d0bc9ae823a80459b46c149e4f51680>, 2025. Notion Blog.
- [16] 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. In *Proceedings of the Twentieth European Conference on Computer Systems*, EuroSys '25, page 1279–1297. ACM, March 2025.
- [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- [18] Weihao Zeng, Yuzhen Huang, Wei Liu, Keqing He, Qian Liu, Zejun Ma, and Junxian He. 7b model and 8k examples: Emerging reasoning with reinforcement learning is both effective and efficient. <https://hkust-nlp.notion.site/simplerl-reason>, 2025. Notion Blog.
- [19] MAA. American invitational mathematics examination - aime. In *American Invitational Mathematics Examination - AIME 2024*, February 2024.
- [20] MAA. American invitational mathematics examination - aime. In *American Invitational Mathematics Examination - AIME 2025*, February 2025.
- [21] MAA. American mathematics competition - amc. In *American Mathematics Competition - AMC*.
- [22] 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.
- [23] 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.
- [24] 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. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems, 2024.
- [25] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps, 2021.
- [26] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022.
- [27] Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. Taco: Topics in algorithmic code generation dataset, 2023.
- [28] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [29] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code, 2024.
- [30] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023.
- [31] Sanghwan Bae, Jiwoo Hong, Min Young Lee, Hanbyul Kim, JeongYeon Nam, and Donghyun Kwak. Online difficulty filtering for reasoning oriented reinforcement learning, 2025.
- [32] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [33] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020.