

SCALER: Synthetic Scalable Adaptive Learning Environment for Reasoning

Anonymous ACL submission

Abstract

Reinforcement learning (RL) offers a principled way to enhance the reasoning capabilities of large language models, yet its effectiveness hinges on training signals that remain informative as models evolve. In practice, RL progress often slows when task difficulty becomes poorly aligned with model capability or when training is dominated by a narrow set of recurring problem patterns. To jointly address these issues, we propose **SCALER** (Synthetic sCalable Adaptive Learning Environment for Reasoning), a framework that sustains effective learning signals through adaptive environment design. SCALER introduces a scalable synthesis pipeline that converts real-world programming problems into verifiable reasoning environments with controllable difficulty and unbounded instance generation, enabling RL training beyond finite datasets while preserving strong correctness guarantees. Building on this, SCALER further employs an adaptive multi-environment RL strategy that dynamically adjusts instance difficulty and curates the active set of environments to track the model’s capability frontier and maintain distributional diversity. This co-adaptation prevents reward sparsity, mitigates overfitting to narrow task patterns, and supports sustained improvement throughout training. Extensive experiments show that SCALER consistently outperforms dataset-based RL baselines across diverse reasoning benchmarks and exhibits more stable, long-horizon training dynamics.

1 Introduction

Reinforcement learning (RL) has become a key post-training paradigm for enhancing the reasoning capabilities of large language models (LLMs) (OpenAI et al., 2024; DeepSeek-AI et al., 2025). By optimizing a policy under explicit, verifiable rewards, RL can sharpen decision making, improve long-horizon credit assignment, and expand the model’s reasoning frontier (Shao et al., 2024; Yu

et al., 2025; Hu, 2025). However, scaling RL for reasoning is often bottlenecked not by the optimizer itself, but by the availability of continuously effective reward signals throughout training (Razin et al., 2023, 2025; Zhang et al., 2025).

In this work, we argue that the training signal, for RL to continue improving LLMs, should remain effective in two complementary senses. First, problems should stay near the model’s current capability boundary during training (Parashar et al., 2025b; Chen et al., 2025d), i.e., remain neither trivial nor unsolvable. When the model mostly sees easy problems, learning saturates; when it mostly sees overly difficult problems, exploration becomes unproductive and the reward becomes sparse. Second, the training distribution should retain sufficient diversity over time (Li et al., 2025). Even if difficulty is well matched, repeatedly interacting with a narrow task distribution can lead to overfitting to a limited set of patterns, weakening generalization and reducing exploration. Note that difficulty variation is not equivalent to diversity. Even if a single environment can generate infinitely many problems by scaling parameters (e.g., longer arrays or larger graphs), it still shares a limited set of templates and failure modes, so learning can plateau once those patterns are mastered.

To this end, we propose **SCALER** (Synthetic sCalable Adaptive Learning Environment for Reasoning), a system that combines scalable reasoning environment synthesis with adaptive multi-environment reinforcement learning. For the first limitation, SCALER provides a way to generate verifiable tasks at scale with controllable difficulty. We develop a synthesis pipeline that programmatically converts real-world programming problems into reasoning environments with (i) verifiable interaction via deterministic oracles and unit tests, (ii) controllable difficulty via explicit scale parameters, and (iii) unbounded instance generation within each environment through randomized testcases

085 generation. The pipeline automatically extracts
086 meta-information, validates testcase generators via
087 breadth/depth checks. This enables scaling training
088 beyond finite reasoning datasets or a small set of
089 hand-crafted environments, while retaining strong
090 correctness guarantees.

091 On the other hand, SCALER designs an adaptive
092 multi-environment RL framework to adapt both in-
093 stance difficulty and environment selection so that
094 training continues to encounter informative chal-
095 lenges as the model improves. Specifically, within
096 each environment, an online difficulty controller
097 adjusts scale parameters based on on-policy rollout
098 accuracy to keep sampled instances near a target
099 success rate, thereby tracking the model’s capabil-
100 ity frontier and avoiding degenerate “all-correct” or
101 “all-wrong” regimes. Across environments, an en-
102 vironment curation mechanism maintains an active
103 set of environments and replaces those whose learn-
104 ing signal has saturated (e.g., difficulty no longer
105 increases or the environment becomes consistently
106 trivial/unlearnable). This realizes an environment-
107 level effectiveness. As the model and environments
108 co-evolve, the marginal learning benefit can dimin-
109 ish, so continuously refreshing the active set helps
110 preserve novelty and sustained learning signals.

111 Our contributions can be summarized as follows:

- 112 • We highlight the importance of difficulty and
113 diversity controllable environments for scaling
114 RL post-training in reasoning.
- 115 • We propose **SCALER** that combines verifi-
116 able, difficulty-controllable environment synthe-
117 sis with adaptive multi-environment RL.
- 118 • Extensive experiments demonstrate SCALER
119 yields consistent improvements across diverse
120 reasoning benchmarks and exhibits more sus-
121 tained training dynamics than dataset-based RL
122 baselines under comparable budgets.

123 2 Related Work

124 **Data-centric RL.** A data-centric line of work im-
125 proves the scalability of RL for LLM reasoning
126 by continuously expanding the training distribu-
127 tion with synthetic data (Wu et al., 2025; Setlur
128 et al., 2024) and self-play (Fang et al., 2025; Liang
129 et al., 2025; Wang et al., 2025b; Chen et al., 2025a).
130 On the synthetic-data side, works like Synthetic
131 Data RL (Guo et al., 2025b) and SWiRL (Goldie
132 et al., 2025) generate task-specific supervision from
133 question–answer pairs to multi-step reasoning and
134 tool-use trajectories, and then apply RL on the re-

135 sulting synthetic corpus. On the self-play side,
136 recent methods produce interactions and hard cases
137 via adversarial games or role-based play, enabling
138 continual data self-generation for policy improve-
139 ment (Zhao et al., 2025; Liu et al., 2025a). How-
140 ever, these approaches are easy to encounter bottle-
141 necks: as the max difficulty of tasks is bounded by
142 the generator agent (Chae et al., 2025), the data may
143 drift out of sync with the evolving policy, collaps-
144 ing into instances that are too easy or too hard and
145 thus weakening learning signals in later stages. In
146 contrast, SCALER applies difficulty controller and
147 environment curation mechanism on synthesized
148 verifiable reasoning environment to keep training
149 informative throughout RL.

Difficulty-aware RL. To mitigate vanishing
150 learning signals, difficulty-aware RL often employs
151 curriculum learning (Team, 2025a; Shi et al., 2025;
152 Chen et al., 2025b) or difficulty scheduling (Wang
153 et al., 2025a; Chen et al., 2025c). Curriculum-
154 based method trains models from easy to hard (Liu
155 et al., 2025b; Parashar et al., 2025a), keeping ef-
156 fective learning signals in training process. How-
157 ever, curricula designed over static dataset faces
158 the challenge of being coarse-grained and hard to
159 design. This motivates environment-based formu-
160 lations where difficulty of each environment can be
161 monitored and adjusted online (Guo et al., 2025a).
162 Reasoning Gym (Stojanovski et al., 2025) offers
163 a suite of procedurally generated, verifiable rea-
164 soning environments with tunable difficulty, while
165 RLVE (Zeng et al., 2025b) further adapts the diffi-
166 culty distribution online as the policy improves.
167 Despite enabling difficulty tracking and adapta-
168 tion, existing environment suites are largely hand-
169 engineered and limited in environment diversity
170 at scale. SCALER addresses this by automatically
171 synthesizing a large and diverse set of reasoning en-
172 vironments with verifiable oracles and controllable
173 difficulty, supporting sustained multi-environment
174 RL.

175 3 SCALER

176 To enable the model to explore, adapt, and general-
177 ize in multiple dynamic environments, we design
178 **SCALER** with two components: (i) an adaptive
179 multi-environment training framework that learns
180 through interaction with a set of verifiable envi-
181 ronments, and (ii) a systematic synthesis pipeline
182 that converts programming problems into difficulty-
183 controllable reasoning environments. Specifically,
184

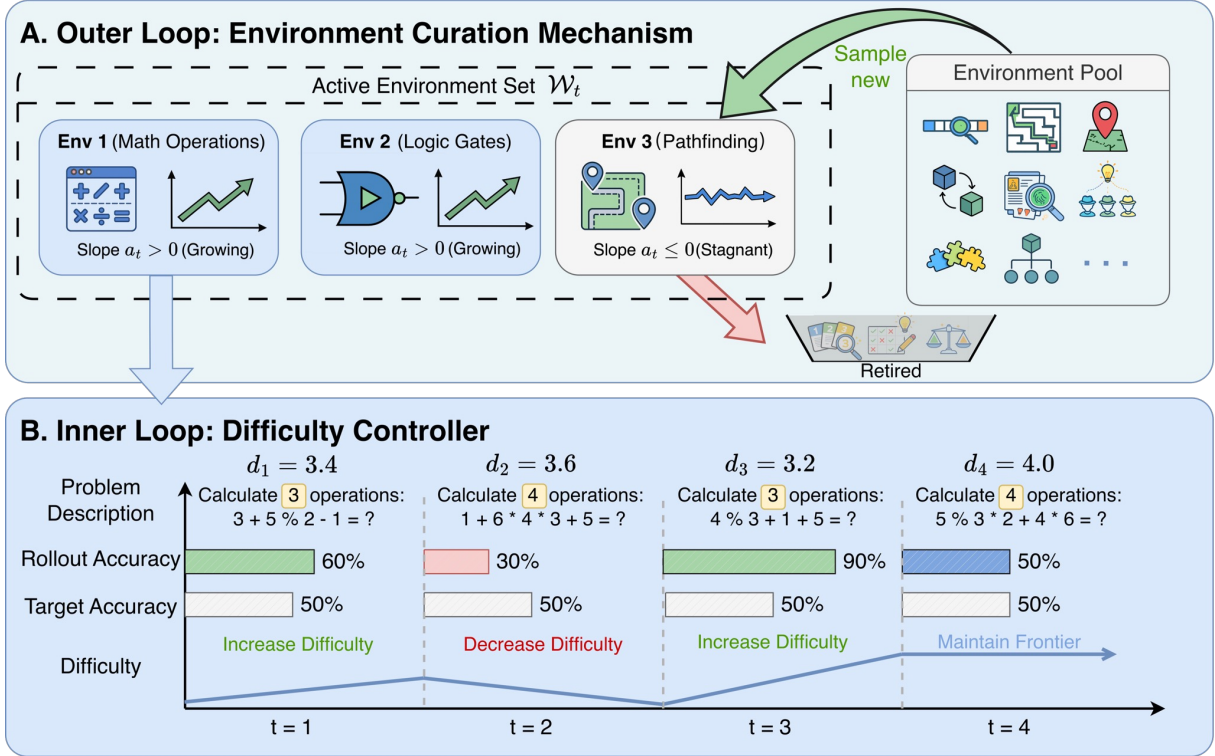


Figure 1: An illustrative training example of **SCALER**. The model interacts with a set of active environments: instances are used for training, online accuracy updates the in-environment difficulty controller, and environments whose learning signal saturates are retired and replaced by new environments via the curation mechanism.

we first introduce the learning strategy over environments, including in-environment difficulty controller and the environment curation mechanism (§3.1), then describe the environment synthesis process from programming problems (§3.2).

3.1 Multi-environment Training Framework

We first illustrate the SCALER environment to facilitate understanding. We define an environment as a context-conditioned, parameterized problem space in which individual instances can evolve through controllable parameters. For example, as shown in the math-operations environment in Figure 1, all problems share a common contextual description, *calculate N operations*, while problem difficulty is modulated by the scale parameter N . More illustrative cases are provided in Appendix B.

To fully unleash the potential of these difficulty-controllable reasoning environments, we adopt two complementary learning strategies: (i) Difficulty Controller. In environments with explicit verifiability, we dynamically adjust the environment’s scale parameters according to the online accuracy, which provides continuous and stable learning signals. (ii) Environment Curation Mechanism. For general environments, we introduce environment curation

mechanism within which the model is trained exclusively on the current active set of environments. As the model and the environment gradually co-adapt and the learning signal diminishes in marginal utility, the current environment is retired and replaced with a newly sampled one.

The detailed implementation of SCALER is provided in Algorithm 1.

3.1.1 Difficulty Controller

A key goal of SCALER is to maintain training at the agent’s capability frontier by continuously sampling instances near the boundary of what the current policy can solve. To achieve this, difficulty is dynamically adjusted within each environment based on current accuracy.

Specifically, the difficulty of the instances is characterized by the array length or the number of edges in the graph and we discretize the difficulty into distinct difficulty levels. Let $acc_t \in [0, 1]$ denote the average accuracy over k sampled instances at step t in a given environment, and let $\tau \in [0, 1]$ denote the target accuracy. The continuous difficulty score $d_t \in \mathbb{R}$, initializing $d_0 = 0$, evolves according to:

$$d_{t+1} = \text{clip}(d_t + \beta \cdot (acc_t - \tau), 0, D), \quad (1)$$

Algorithm 1 SCALER Training Framework

Require: Policy π_θ , Environment Pool \mathcal{E}_{pool} , Target accuracy τ , Step size β , Active set size M

```

1: Initialize: Sample initial active set  $\mathcal{W} \leftarrow \{e_1, \dots, e_M\} \subset \mathcal{E}_{pool}$ , Difficulty  $d \leftarrow d_{min}$  for all  $e \in \mathcal{W}$ 
2: for step  $t = 1, 2, \dots, T_{max}$  do
3:    $\mathcal{D}_{batch} \leftarrow \emptyset$ 
4:   for each environment  $e \in \mathcal{W}$  do
5:     Sample instances  $x$  from  $e$  based on difficulty  $d$ 
6:     Collect trajectories using  $\pi_\theta$  and calculate accuracy  $acc_t$ 
7:      $\mathcal{D}_{batch} \leftarrow \mathcal{D}_{batch} \cup \text{Trajectories}$ 
8:     Update difficulty  $d$  for this environment ▷ Eq. 1
9:   Update  $\pi_\theta$  with  $\mathcal{D}_{batch}$ 
10:  for each environment  $e \in \mathcal{W}$  do
11:    Calculate slope  $a_t^e$  over last  $K_{slope}$  steps ▷ Eq. 3
12:    if  $a_t^e \leq 0 \vee acc = 0$  for  $K_{zero}$  steps  $\vee d = D^e$  for  $K_{sat}$  steps then
13:      Retire  $e$ :  $\mathcal{W} \leftarrow \mathcal{W} \setminus \{e\}$ ;  $\mathcal{E}_{pool} \leftarrow \mathcal{E}_{pool} \cup \{e\}$ 
14:      Sample new  $e_{new} \sim \mathcal{E}_{pool}$ 
15:       $\mathcal{W} \leftarrow \mathcal{W} \cup \{e_{new}\}$ 

```

where $\beta > 0$ controls the adaptation rate and $D \in \mathbb{N}$ is the maximum discrete difficulty level in the environment. This update rule increases difficulty when $acc_t > \tau$ and decreases it otherwise, thereby maintaining the training distribution near the agent’s capability boundary.

For the next sampling step, while d_{t+1} is real-valued, the environment can only generate k instances at integer-valued levels in $\{0, 1, \dots, D\}$. To approximate $d_{t+1} \in \mathbb{R}$ under this integer constraint, we construct a multiset of k integers whose mean closely matches d_{t+1} . Specifically, let $\ell = \lfloor d_{t+1} \rfloor$ and define $h = \text{round}(k(d_{t+1} - \ell))$. We assign h instances the value $\ell + 1$ and the remaining $k - h$ instances the value ℓ , thereby balancing rounding up and rounding down; for example, when $d_{t+1} = 2.3$ and $k = 10$, this yields three instances at level 3 and seven at level 2, whose average equals 2.3.

The proposed controller naturally tracks model competence. As the model improves, online accuracy rises above the target and difficulty is pushed upward to stay near the performance boundary. Conversely, if the model forgets or the policy distribution shifts, accuracy drops and the controller reduces difficulty, recovering usable learning signals without manual retuning.

3.1.2 Environment Curation Mechanism

To keep sufficient diversity during training and avoid a "tourist" learning pattern where the model randomly explores a variety of environments without sustained progress, SCALER introduces the

environment curation mechanism which maintains an active set of environments and restricts training to this dynamic set, continuously refreshing it to avoid spending budget on environments that no longer provide meaningful learning signals.

Let $\mathcal{W}_t = \{e_1, \dots, e_{|\mathcal{W}|}\}$ denote the active environment set at training step t . At each step, instances are sampled only from environments in \mathcal{W}_t and used for policy optimization. Environments are periodically evaluated at every step and may be retired from the set and replaced by newly synthesized environments.

For each environment e , let d_t^e denote its difficulty at step t . To detect stalled progress, we argue that when the recent difficulty-step slope is less or equal to zero, the learning on this environment is saturable and should be retired. Specifically, we estimate the slope by fitting a least-squares line over the last K_{slope} steps, defining the index set $T = \{t - K_{slope} + 1, \dots, t\}$ and fitting

$$d_i^e \approx a_t^e \cdot i + b_t^e, \quad i \in T, \quad (2)$$

where the slope a_t^e is given in closed form by

$$a_t^e = \frac{\sum_{i \in T} (i - \bar{i})(d_i^e - \bar{d}^e)}{\sum_{i \in T} (i - \bar{i})^2}, \quad (3)$$

with $\bar{i} = \frac{1}{|T|} \sum_{i \in T} i$ and $\bar{d}^e = \frac{1}{|T|} \sum_{i \in T} d_i^e$. An environment is retired if $a_t^e \leq 0$, indicating that difficulty is temporarily saturated in that environment.

To accelerate retirement decisions, two additional heuristics are applied. Let $acc_t^e \in [0, 1]$

denote the accuracy at step t and let D^e be the maximum difficulty of environment e . An environment is retired early if it is unlearnable, i.e., $\text{acc}^e = 0$ for K_{zero} consecutive steps, or saturated, i.e., $d^e = D^e$ for K_{sat} consecutive steps.

When an environment is retired, it is replaced by a newly synthesized environment sampled from the environment pool, keeping the set size fixed. Such retirement is temporary, as it targets situations where the model has reached a plateau in learning within the current environment, rendering further progress unproductive. However, as training continues and the model’s capability boundary improve, there remains potential for continued learning within these environments. Therefore, retired environments are reintegrated into the environment pool and will be resampled in future iterations, offering renewed opportunities for the model to explore and adapt.

In general, the environment curation mechanism supports sustained learning within multiple environments while preventing training from stagnating on environments that have become either trivial or unproductive.

3.2 From Code Generation to Reasoning Environments

Our multi-environment training framework continuously tracks the model’s capability frontier while maintaining data diversity and freshness. This requires a steady supply of environments that are difficulty-controllable and automatically verifiable. We therefore introduce an environment synthesis pipeline that converts real-world programming problems into such environments.

Accordingly, our synthesis pipeline is organized into three components:

- **Meta information extraction** (§3.2.1). Ensures rewardable and unambiguous supervision by selecting problems with well-defined, verifiable outputs, mitigating output-format mismatch and reward hacking.
- **Testcase generation and verification** (§3.2.2). Provides a stream of diverse, valid, and automatically verifiable instances, addressing the mismatch between static problem statements and stochastic environment inputs.
- **Heuristic difficulty calibration** (§3.2.3). Defines an executable difficulty range under context and runtime budgets, addressing the scale mismatch between real-world inputs and agentic prompting.

3.2.1 Extracting Meta Information

To offer high-quality programming problems for reasoning environment synthesis and mitigate the risk of reward hacking, we adopt a prompt-based method to extract key meta-information from candidate problems, followed by a rule-based filtering process. The complete extraction prompt is provided in Appendix A.1.

For each programming problem p , we extract a metadata tuple $\text{meta}(p)$ comprising (i) scale parameters that characterize complexity-related constraint variables, and (ii) output requirements that describe the output type and whether the correct output is unique. Accordingly, we discard ill-formed problems and retain only those with a unique output whose type lies in $\{\text{number}, \text{array}, \text{string}\}$.

3.2.2 Generating and Verifying Testcases

To sample different instances with fixed extracted scale parameters from one specific environment, we construct a testcase generator agent that takes a target scale parameter configuration and produces an input whose content is randomized but conforms to the original problem specification. We provide the full prompting details in Appendix A.2. Since the environment relies on the synthetic agent, we validate generator functions with two complementary checks.

Breadth check. We sample diverse scale parameters and corresponding inputs, when following (Fu et al., 2025), each instance is evaluated by multiple independent ground-truth solutions. Consistency across solutions is required, which simultaneously (i) detects malformed generator function and (ii) revalidates the assumption that the output is unique.

Deep check. For a fixed scale parameter configuration, we call generator function multiple times and compare the resulting ground-truth outputs across generated instances. Based on the numbers and max size requirement of output clustering, we enforce sufficient diversity to prevent the training process from overfitting to a narrow pattern distribution or exploiting reward hacking.

3.2.3 Calibrating Heuristic Difficulty Levels

As inputs for real-world programming problems can be millions, which is not suitable for the prompt used in agentic training, we aim to re-define the difficulty range for each environment that is both scalable and executable. Apparently, the effective

394 maximum difficulty is bounded by two practical
395 constraints: (i) the maximum prompt length ac-
396 cepted by the policy model, and (ii) the execution
397 time limit enforced by the original problem setting.

398 To estimate the largest feasible scale parameter
399 configuration s_{\max} under the prompt-length con-
400 straint, a binary search is performed over a single
401 global scale factor. Feasibility is determined by
402 whether the testcase fits within the context window
403 while respecting the fixed execution time limit.

404 Given s_{\min} and s_{\max} , the environment defines
405 a finite set of difficulty levels by discretizing the
406 scaling range. To keep the number of levels within
407 a reasonable budget, the discretization strategy de-
408 pends on the span of each parameter: when the
409 range is small, levels follow an arithmetic progres-
410 sion; when the range is large, levels follow a geo-
411 metric progression. This heuristic yields a compact
412 but expressive difficulty ladder, enabling smooth
413 scheduling while preserving meaningful granular-
414 ity across scales.

415 4 Experiments

416 SCALER addresses two key challenges: (i) pro-
417 viding a synthesis pipeline that systematically
418 generates environments with infinite data gener-
419 ation capabilities, and (ii) developing a multi-
420 environment training framework, which enables
421 sustaining model improvement through adaptive
422 learning across diverse environments. To valid the
423 efficiency of our approach, experiments are orga-
424 nized around three research questions (RQs):

- 425 • RQ1: How does SCALER compare with dataset-
426 based training baselines under comparable train-
427 ing budgets?
- 428 • RQ2: How does performance vary as the number
429 of environments increases?
- 430 • RQ3: Are all components of SCALER Necessary
431 for their gains?

432 4.1 Experimental Setup

433 **Synthesis Pipeline Settings.** We use CodeCon-
434 tests (Li et al., 2022) as the seed dataset and
435 GLM-4.6 (Zeng et al., 2025a) as the data-synthesis
436 agent. For code execution, we use Sandbox-
437 Fusion (Bytedance-Seed-Foundation-Code-Team
438 et al., 2025). After filter, we obtain the subset
439 consisting of 4973 programming problems and syn-
440 thesize 2739 SCALER environments. Specifically,
441 we offer environment cases in Appendix B.

442 **Training setup.** All experiments start from
443 Qwen3 (Team, 2025b) series, and we use Qwen3-
444 1.7B-base and Qwen3-4B-base as the policy model.
445 For the environment curation mechanism, we set
446 $K_{slope} = 10$ and $K_{zero} = K_{sat} = 5$. The train-
447 ing batch size equals the environment set size: at
448 each optimization step, one problem is sampled
449 from each of the 64 environments according to
450 the environment-specific difficulty controller, re-
451 sulting in 64 prompts per step. Reinforcement
452 learning is performed with GRPO (Shao et al.,
453 2024). For each prompt, $n_{resp}=8$ responses are
454 sampled. Training-time decoding uses tempera-
455 ture $T=1.0$ and $top_p=1.0$. The prompt length
456 budget is capped at 4096 tokens and the response
457 length budget at 8192 tokens. Additional hyperpa-
458 rameters and infrastructure details are deferred to
459 Appendix C.

460 **Evaluation protocol.** Performance is evaluated
461 on five benchmarks: AIME24 (Zhang and Math-
462 AI, 2024), AMC23 (Art of Problem Solving, 2023),
463 MATH-500 (Lightman et al., 2023), MMLU-
464 Pro (Wang et al., 2024), and BBEH (Kazemi
465 et al., 2025). Results are reported as avg@16 for
466 AIME24, AMC23 and MMLU-Pro, avg@1 for the
467 remaining benchmarks. Unless otherwise spec-
468 ified, decoding uses temperature $T = 0.6$ and
469 $top_p = 0.95$.

470 4.2 RQ1: Environment-based training is 471 better than Dataset-based training

472 A common paradigm for improving reasoning is to
473 scale up training on static datasets, such as curated
474 corpora with verified solutions or collections aug-
475 mented by teacher-generated answers. However,
476 static corpora provide limited effective learning sig-
477 nals, where agents must repeatedly adapt to new
478 environments and generalize beyond previously
479 seen distributions. To evaluate whether interactive
480 environment-based learning better supports the im-
481 provement of model capability, this RQ compares
482 SCALER with two dataset-based baselines under
483 comparable training budgets: MATH (Hendrycks
484 et al., 2021), a curated corpus with verified solu-
485 tions, and DeepMath (He et al., 2025), a dataset of
486 teacher-generated solutions.

487 Table 1 shows that both dataset baselines sub-
488 stantially improve over the Qwen3-1.7B-Base and
489 Qwen3-4B-Base, with particularly large gains on
490 math-centric benchmarks. In contrast, SCALER
491 achieves the best overall average and delivers con-

Model	MATH-500	AMC23	AIME24	MMLU-Pro	BBEH	AVG
Qwen3-1.7B-base	59.6	29.21	3.33	33.30	3.26	25.74
+DeepMath-103K	73.6	47.97	14.58	49.64	9.56	39.07
+MATH-7.5k	75.6	50.78	15.20	46.78	6.08	38.89
+SCALER	75.8	49.53	12.91	50.89	11.74	40.18
Qwen3-4B-Base	66.4	44.70	8.75	51.60	8.10	35.91
+DeepMath-103K	86.6	65.60	22.29	68.03	12.82	51.08
+MATH-7.5k	86.2	70.63	24.16	69.19	10.00	52.04
+SCALER	84.4	75.0	27.29	70.00	14.56	54.25

Table 1: Performance comparison of Dataset-based RL and SCALER on five reasoning benchmarks. AVG is the unweighted mean over the five benchmarks. The highest performance is **bolded**.

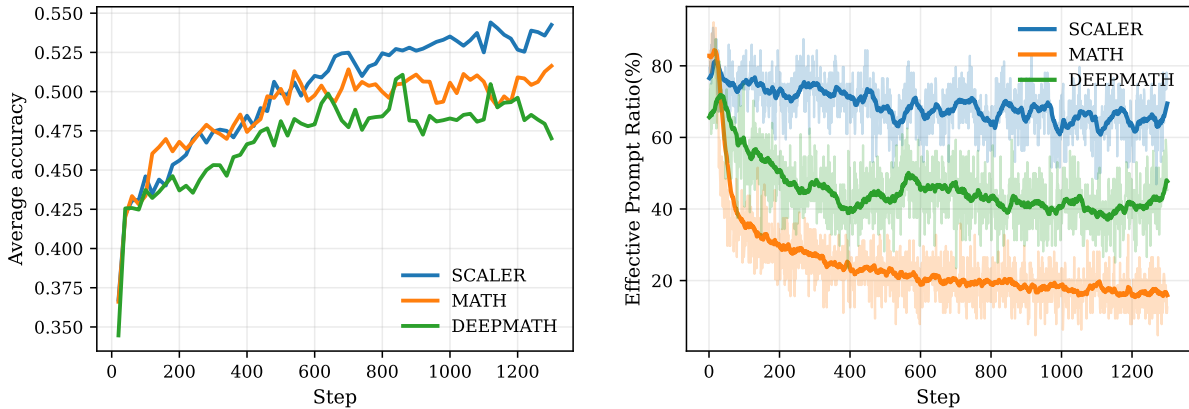


Figure 2: Left: average performance across the five evaluation benchmarks during Qwen3-4B-base training, comparing dataset-based baselines (MATH, DeepMath) and SCALER. Right: effective sampling statistics under SCALER, indicating that most sampled instances remain near the model’s capability boundary.

492 sistent improvements across all five evaluations,
493 including three math benchmarks, MMLU-Pro and
494 BBEH, suggesting stronger transfer beyond nar-
495 rowly curated math data.

496 The training dynamics in Figure 2 provide addi-
497 tional evidence for this gap. Figure 2 (left) shows
498 that SCALER not only reaches a higher level, but
499 also exhibits a more sustained improvement trend:
500 the aggregated evaluation performance continues
501 to rise for more than 1,000 training steps, whereas
502 dataset-based baselines plateau earlier. Importantly,
503 this behavior is supported by boundary-focused
504 sampling: Figure 2 (right) indicates that SCALER
505 keeps a higher effective sampling rate than dataset-
506 based baselines, where most prompts remain nei-
507 ther trivial nor intractable. By continuously provid-
508 ing instances near the model’s current capability
509 boundary from SCALER environments, SCALER
510 sustains informative reward during RL, mitigates
511 premature saturation, and enables longer-horizon
512 performance gains.

4.3 RQ2: Scaling environment size leads to incremental performance gains

513 SCALER introduces the automatic environment
514 synthesis pipeline which significantly reduces labor
515 costs. This RQ examines how model performance
516 changes as the number of environments increases,
517 with results summarized in Figure 3.
518

519 We conduct the experiment by randomly sample
520 8, 64, 512 environments from 2739 SCALER en-
521 vironments, where the larger size of environments
522 always contain the smaller size one. As shown in
523 Figure 3, increasing the number of environments
524 from 8 to 2739 leads to incremental performance
525 improvements. The model benefits from encoun-
526 tering a greater diversity of tasks, which allows
527 it to maintain consistent learning progress. With
528 exposure to a wider variety of challenges and en-
529 vironments facilitating the development of broader
530 reasoning skills, scaling the number of environ-
531 ments enables the model to continuously adapt and
532 enhance its reasoning capabilities.
533

534 It is worth noting that even with a smaller num-

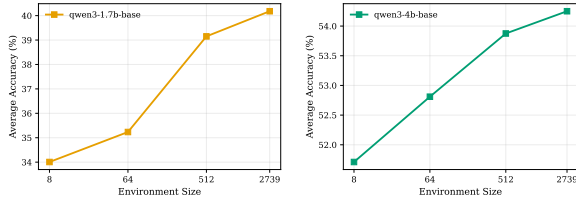


Figure 3: Accuracy improvements for both Qwen3-4B-base and Qwen3-1.7B-base as environment size increases. Both models show a consistent increase in performance with larger environment sizes.

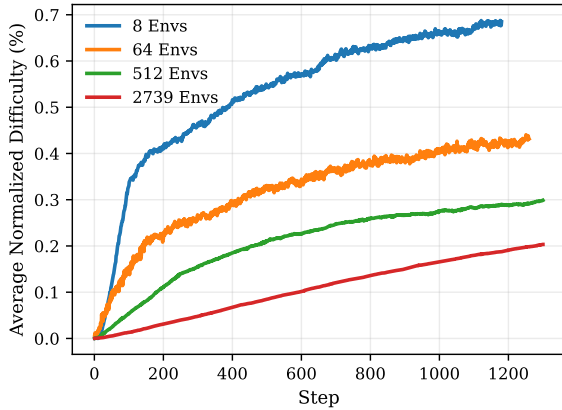


Figure 4: Training dynamics of Qwen3-4B-Base with different numbers of environments. The plot shows that even with smaller environments, the model continues to learn with increasing difficulty levels.

535 ber of environments, the model continues to engage
 536 in ongoing exploration, with increasing difficulty
 537 levels within each individual environment, as illus-
 538 trated in Figure 4. The key insight is that scaling
 539 the number of environments with fixed training
 540 budgets, inherently involves a trade-off between
 541 task difficulty and diversity. Striking right balance
 542 is critical: focusing solely on increasing difficulty
 543 may result in limited improvements, while exces-
 544 sive diversity without a proper difficulty controller,
 545 i.e. DeepMath dataset with 103k samples, does not
 546 necessarily lead to optimal performance.

547 4.4 RQ3: All components of SCALER are 548 necessary for its gains

549 This experiment studies whether SCALER’s im-
 550 provements require its two core components in
 551 multi-environment training framework: adaptive
 552 difficulty controller and environment curation
 553 mechanism.

554 Two ablations based on Qwen3-4B-base are
 555 compared against the full SCALER system in Fig-
 556 ure 5. Removing the curation mechanism disables

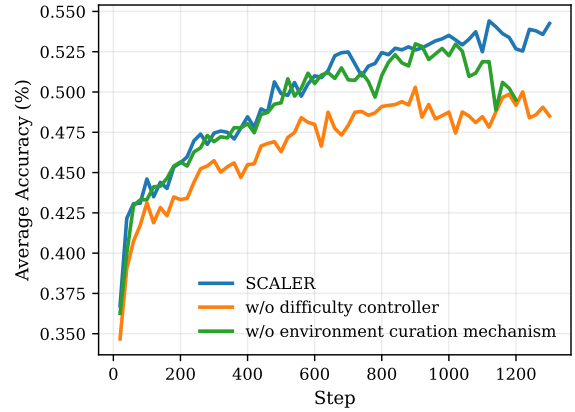


Figure 5: Ablation study of SCALER on Qwen3-4B-Base. Validation accuracy over training steps for the full system and two variants that remove difficulty controller or the environment curation mechanism. Both components contribute to stronger and more sustained performance improvements.

environment replacement, forcing training to ran- 557
 558 dom sample every environment. To model the abla-
 559 tion of difficulty controller, we replace the bound-
 560 ary tracking with random difficulty sampling, se-
 561 lecting 5 instances per environment. This modifi-
 562 cation reduces the consistency of the supervision
 563 signal, as the environment selection is no longer
 564 guided by a difficulty-based approach.

565 Both ablations lead to lower performance than
 566 the full system, indicating that difficulty controller
 567 and the curation mechanism each make a material
 568 contribution. Qualitatively, difficulty controller pre-
 569 vents training from drifting into regimes that are
 570 too easy or too hard, while the curation mecha-
 571 nism promotes sustained learning within learnable
 572 environments and mitigates shallow adaptation.

573 5 Conclusion

574 In this paper, we proposes a novel approach
 575 SCALER containing two key components: (i) rea-
 576 soning environment synthesis pipeline, (ii) multi-
 577 environment training framework. Through our syn-
 578 thesis pipeline to generating a variety of difficulty-
 579 controllable environments , our approach offers a
 580 platform for research community to explore the im-
 581 pact of properties of environments on RL training.
 582 Moreover, the multi-environment training frame-
 583 work ensures the difficulty of instances consistently
 584 matches the model’s capabilities and maintains di-
 585 versity and freshness, jointly preserving the **contin-**
 586 **uously effective** reward signals to improve model
 587 capabilities in long horizon.

588 Limitations

589 While our work provides a comprehensive anal-
590 ysis of SCALER’s ability to improve model per-
591 formance through dynamic environment synthesis
592 and multi-environment training, there are several
593 limitations that warrant further exploration:

- 594 • **Exploration of Environment Internal At-**
595 **tributes:** Our study primarily focuses on scal-
596 ing environment size and its components. How-
597 ever, the internal characteristics of environments,
598 such as the richness of context, intrinsic diffi-
599 culty, and other environment-specific properties,
600 have yet to be thoroughly investigated. Future
601 research should examine how these factors influ-
602 ence model performance, especially in the con-
603 text of dynamic difficulty adjustment and envi-
604 ronment transitions during training.
- 605 • **Limited Scope of Environments:** Our study
606 is based on 2739 SCALER environments, and
607 while this is a larger number compared to hand-
608 crafted ones, the scaling still hasn’t subsided due
609 to our experiments. Further research is needed to
610 explore scaling laws related to environment size,
611 model size, and computational resources, and to
612 understand how these scaling factors impact both
613 model training efficiency and performance.

614 References

615 Art of Problem Solving. 2023. Amc 12 (2023)
616 problems. [https://artofproblemsolving.com/
617 wiki/index.php/2023_AMC_12A_Problems](https://artofproblemsolving.com/wiki/index.php/2023_AMC_12A_Problems). We
618 use the 2023 AMC 12A/12B problems as hosted
619 by AoPS; a processed version is available at [https:
620 //huggingface.co/datasets/zwhe99/amc23](https://huggingface.co/datasets/zwhe99/amc23).

621 Bytedance-Seed-Foundation-Code-Team, :, Yao Cheng,
622 Jianfeng Chen, Jie Chen, Li Chen, Liyu Chen, Wen-
623 tao Chen, Zhengyu Chen, Shijie Geng, Aoyan Li,
624 Bo Li, Bowen Li, Linyi Li, Boyi Liu, Jiaheng Liu,
625 Kaibo Liu, Qi Liu, Shukai Liu, and 37 others. 2025.
626 *Fullstack bench: Evaluating llms as full stack coders*.
627 *Preprint*, arXiv:2412.00535.

628 Justin Yang Chae, Md Tanvirul Alam, and Nidhi Ras-
629 togi. 2025. *Towards understanding self-play for llm
630 reasoning*. *Preprint*, arXiv:2510.27072.

631 Jiaqi Chen, Bang Zhang, Ruotian Ma, Peisong Wang,
632 Xiaodan Liang, Zhaopeng Tu, Xiaolong Li, and
633 Kwan-Yee K. Wong. 2025a. *Spc: Evolving self-
634 play critic via adversarial games for llm reasoning*.
635 *Preprint*, arXiv:2504.19162.

636 Xiaoyin Chen, Jiarui Lu, Minsu Kim, Dinghui Zhang,
637 Jian Tang, Alexandre Piché, Nicolas Gontier, Yoshua
638 Bengio, and Ehsan Kamaloo. 2025b. *Self-evolving*

curriculum for llm reasoning. *arXiv preprint*
arXiv:2505.14970. 639 640

Xinjie Chen, Minpeng Liao, Guoxin Chen, Chengxi
Li, Biao Fu, Kai Fan, and Xinggao Liu. 2025c.
*From data-centric to sample-centric: Enhancing llm
reasoning via progressive optimization*. *Preprint*,
arXiv:2507.06573. 641 642 643 644 645

Zhuoyue Chen, Jihai Zhang, Ben Liu, Fangquan Lin,
and Wotao Yin. 2025d. *Scale down to speed up:
Dynamic data selection for reinforcement learning*.
Training, 2500:3000. 646 647 648 649

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.
2025. *Deepseek-r1: Incentivizing reasoning capa-
bility in llms via reinforcement learning*. *Preprint*,
arXiv:2501.12948. 650 651 652 653 654 655 656 657

Wenkai Fang, Shunyu Liu, Yang Zhou, Kongcheng
Zhang, Tongya Zheng, Kaixuan Chen, Mingli Song,
and Dacheng Tao. 2025. *Serl: Self-play reinforce-
ment learning for large language models with limited
data*. *Preprint*, arXiv:2505.20347. 658 659 660 661 662

Jia Fu, Xinyu Yang, Hongzhi Zhang, Yahui Liu,
Jingyuan Zhang, Qi Wang, Fuzheng Zhang, and
Guorui Zhou. 2025. *Klear-codetest: Scalable test
case generation for code reinforcement learning*.
Preprint, arXiv:2508.05710. 663 664 665 666 667

Anna Goldie, Azalia Mirhoseini, Hao Zhou, Irene Cai,
and Christopher D. Manning. 2025. *Synthetic data
generation & multi-step rl for reasoning & tool use*.
Preprint, arXiv:2504.04736. 668 669 670 671

Jiacheng Guo, Ling Yang, Peter Chen, Qixin Xiao, Yin-
jie Wang, Xinzhe Juan, Jiahao Qiu, Ke Shen, and
Mengdi Wang. 2025a. *Genenv: Difficulty-aligned
co-evolution between llm agents and environment
simulators*. *Preprint*, arXiv:2512.19682. 672 673 674 675 676

Yiduo Guo, Zhen Guo, Chuanwei Huang, Zi-Ang
Wang, Zekai Zhang, Haofei Yu, Huishuai Zhang, and
Yikang Shen. 2025b. *Synthetic data rl: Task defini-
tion is all you need*. *Preprint*, arXiv:2505.17063. 677 678 679 680

Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu
Chen, Yue Wang, Linfeng Song, Dian Yu, Zhen-
wen Liang, Wenxuan Wang, and 1 others. 2025.
*Deepmath-103k: A large-scale, challenging, decon-
taminated, and verifiable mathematical dataset for ad-
vancing reasoning*. *arXiv preprint arXiv:2504.11456*. 681 682 683 684 685 686

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021. *Measuring mathematical
problem solving with the math dataset*. *Preprint*,
arXiv:2103.03874. 687 688 689 690 691

Jian Hu. 2025. *Reinforce++: A simple and efficient
approach for aligning large language models*. *arXiv
preprint arXiv:2501.03262*. 692 693 694

695	Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, Sanket Vaibhav Mehta, Lalit K Jain, Virginia Aglietti, Disha Jindal, Peter Chen, Nishanth Dikkala, Gladys Tyen, Xin Liu, Uri Shalit, Silvia Chiappa, Kate Olszewska, Yi Tay, Vinh Q. Tran, Quoc V. Le, and Orhan Firat. 2025. BIG-bench extra hard . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 26473–26501, Vienna, Austria. Association for Computational Linguistics.	James Caverlee, Dileep Kalathil, and Shuiwang Ji. 2025a. Curriculum reinforcement learning from easy to hard tasks improves llm reasoning . <i>Preprint</i> , arXiv:2506.06632.	753 754 755 756
696		Shubham Parashar, Shurui Gui, Xiner Li, Hongyi Ling, Sushil Vemuri, Blake Olson, Eric Li, Yu Zhang, James Caverlee, Dileep Kalathil, and 1 others. 2025b. Curriculum reinforcement learning from easy to hard tasks improves llm reasoning. <i>arXiv preprint arXiv:2506.06632</i> .	757 758 759 760 761 762
697		Noam Razin, Zixuan Wang, Hubert Strauss, Stanley Wei, Jason D Lee, and Sanjeev Arora. 2025. What makes a reward model a good teacher? an optimization perspective. <i>arXiv preprint arXiv:2503.15477</i> .	763 764 765 766
698		Noam Razin, Hattie Zhou, Omid Saremi, Vimal Thilak, Arwen Bradley, Preetum Nakkiran, Joshua Susskind, and Etai Littwin. 2023. Vanishing gradients in reinforcement finetuning of language models. <i>arXiv preprint arXiv:2310.20703</i> .	767 768 769 770 771
699		Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. 2024. RL on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold . <i>Preprint</i> , arXiv:2406.14532.	772 773 774 775 776
700		Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	777 778 779 780 781 782
701		Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. 2025. Efficient reinforcement finetuning via adaptive curriculum learning. <i>arXiv preprint arXiv:2504.05520</i> .	783 784 785 786
702		Zafir Stojanovski, Oliver Stanley, Joe Sharratt, Richard Jones, Abdulhakeem Adefioye, Jean Kaddour, and Andreas Köpf. 2025. Reasoning gym: Reasoning environments for reinforcement learning with verifiable rewards. <i>arXiv preprint arXiv:2505.24760</i> .	787 788 789 790 791
703		Kimi Team. 2025a. Kimi k1.5: Scaling reinforcement learning with llms.	792 793
704		Qwen Team. 2025b. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.	794 795
705		Hong Wang, Zhezheng Hao, Jian Luo, Chenxing Wei, Yao Shu, Lei Liu, Qiang Lin, Hande Dong, and Jiawei Chen. 2025a. Scheduling your llm reinforcement learning with reasoning trees . <i>Preprint</i> , arXiv:2510.24832.	796 797 798 799 800
706	Shipeng Li, Shikun Li, Zhiqin Yang, Xinghua Zhang, Gaode Chen, Xiaobo Xia, Hengyu Liu, and Zhe Peng. 2025. Learnalign: Reasoning data selection for reinforcement learning in large language models based on improved gradient alignment . <i>Preprint</i> , arXiv:2506.11480.	Pinzheng Wang, Juntao Li, Zecheng Tang, Haijia Gui, and Min zhang. 2025b. Improving rationality in the reasoning process of language models through self-playing game . <i>Preprint</i> , arXiv:2506.22920.	801 802 803 804
707			
708			
709			
710			
711			
712	Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Axel Gimeno Gil, 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, and 7 others. 2022. Alphacode data materials . Includes dm-code_contests (CodeContests competitive programming dataset used in AlphaCode).		
713			
714			
715			
716			
717			
718			
719			
720			
721			
722	Xiao Liang, Zhong-Zhi Li, Yeyun Gong, Yang Wang, Hengyuan Zhang, Yelong Shen, Ying Nian Wu, and Weizhu Chen. 2025. Sws: Self-aware weakness-driven problem synthesis in reinforcement learning for llm reasoning . <i>Preprint</i> , arXiv:2506.08989.		
723			
724			
725			
726			
727	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step . <i>arXiv preprint arXiv:2305.20050</i> . Releases PRM800K and the MATH-500 (500-problem) test subset via https://github.com/openai/prm800k .		
728			
729			
730			
731			
732			
733			
734	Bo Liu, Leon Guertler, Simon Yu, Zichen Liu, Penghui Qi, Daniel Balcells, Mickel Liu, Cheston Tan, Weiyan Shi, Min Lin, and 1 others. 2025a. Spiral: Self-play on zero-sum games incentivizes reasoning via multi-agent multi-turn reinforcement learning . <i>arXiv preprint arXiv:2506.24119</i> .		
735			
736			
737			
738			
739			
740	Huanyu Liu, Jia Li, Hao Zhu, Kechi Zhang, Yihong Dong, and Ge Li. 2025b. Saturn: Sat-based reinforcement learning to unleash language model reasoning . <i>arXiv preprint arXiv:2505.16368</i> .		
741			
742			
743			
744	OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, and 244 others. 2024. Openai o1 system card . <i>Preprint</i> , arXiv:2412.16720.		
745			
746			
747			
748			
749			
750			
751	Shubham Parashar, Shurui Gui, Xiner Li, Hongyi Ling, Sushil Vemuri, Blake Olson, Eric Li, Yu Zhang,		
752			

805 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,
806 Abhranil Chandra, Shiguang Guo, Weiming Ren,
807 Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max
808 Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang
809 Yue, and Wenhui Chen. 2024. *Mmlu-pro: A more ro-
810 bust and challenging multi-task language understand-
811 ing benchmark*. *arXiv preprint arXiv:2406.01574*.
812 Accepted at NeurIPS 2024 Datasets and Bench-
813 marks (Spotlight). Dataset: [https://huggingface.
814 co/datasets/TIGER-Lab/MMLU-Pro](https://huggingface.co/datasets/TIGER-Lab/MMLU-Pro).

815 Zijian Wu, Jinjie Ni, Xiangyan Liu, Zichen Liu, Hang
816 Yan, and Michael Qizhe Shieh. 2025. *Synthrl: Scal-
817 ing visual reasoning with verifiable data synthesis*.
818 *Preprint*, arXiv:2506.02096.

819 Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,
820 Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,
821 Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin,
822 Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan
823 Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and
824 16 others. 2025. *Dapo: An open-source llm re-
825 inforcement learning system at scale*. *Preprint*,
826 arXiv:2503.14476.

827 Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin
828 Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao
829 Zeng, Jiajie Zhang, and 1 others. 2025a. *Glm-4.5:
830 Agentic, reasoning, and coding (arc) foundation mod-
831 els*. *arXiv preprint arXiv:2508.06471*.

832 Zhiyuan Zeng, Hamish Ivison, Yiping Wang, Lifan
833 Yuan, Shuyue Stella Li, Zhuorui Ye, Siting Li, Jacque-
834 line He, Runlong Zhou, Tong Chen, and 1 others.
835 2025b. *Rlve: Scaling up reinforcement learning for
836 language models with adaptive verifiable environ-
837 ments*. *arXiv preprint arXiv:2511.07317*.

838 Bonan Zhang, Zhongqi Chen, Bowen Song, Qinya
839 Li, Fan Wu, and Guihai Chen. 2025. *Con-
840 fclip: Confidence-weighted and clipped reward
841 for reinforcement learning in llms*. *Preprint*,
842 arXiv:2509.17730.

843 Yifan Zhang and Team Math-AI. 2024. *Aime 2024*.
844 [https://huggingface.co/datasets/math-ai/
845 aime24](https://huggingface.co/datasets/math-ai/aime24).

846 Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin
847 Xu, Matthieu Lin, Shenzi Wang, Qingyun Wu, Zi-
848 long Zheng, and Gao Huang. 2025. *Absolute zero:
849 Reinforced self-play reasoning with zero data*. *arXiv
850 preprint arXiv:2505.03335*.

A Prompt Design

A.1 Extract Meta Information Prompt

Prompt 1: Prompt used for extracting meta information.

```

You are given as input the full statement of a single algorithmic problem. Your job
is to emit one JSON code block that captures:
1. All numeric scale parameters that are relevant to the time complexity of
   typical solutions.
2. The type of the required output.
3. Whether the required output is unique, whether the problem has exactly one
   correct output for each valid input.
## 1. Scale parameters
A scale parameter is any integer quantity that bounds:
- The number of items, elements, or positions (n= number of elements, m = number of
  edges, q = number of queries).
- The size (length) of a grid, string, or sequence (e.g. length up to 2e5, grid up
  to 1000 * 1000).
- The size of a state space or iteration space that an algorithm must explicitly
  handle.
### What to INCLUDE as scale parameters
Include a parameter only if all of the following are true:
- It directly bounds the size or count of something that is iterated over, e.g.:
  - number of elements / vertices / edges / queries (n, m, q, etc.),
  - length of a string or array,
  - rows / columns of a grid.
- The bounds appear explicitly in the statement, usually in the Input / Constraints
  section, like:
  -  $1 \leq n \leq 2 \cdot 10^5$ 
  -  $0 \leq m \leq 2 \cdot 10^5$ 
- You can clearly identify the parameter name (e.g., n, m, q, k, N, etc.).
### What to EXCLUDE from scale parameters
- Number of test cases / groups: never include t, T, or similar when it means "
  number of test cases".
- Pure value ranges for single items that do not change the input size, e.g.:
  -  $-10^9 \leq a_i \leq 10^9$  when  $a_i$  is just the value of an element.
  - Coordinate or weight ranges that are not used as sizes of arrays/grids.
- Any quantity that only affects output format or precision.
### Representation of scale parameters
In the JSON, represent scale parameters under the key "scale_params" as:
<scale_params example>
## 2. Output type classification
You must classify the type of the required output into exactly one of the following
strings:
- "string": The required output is a single string or a small number of strings.
- "number": The required output is a single numeric value (integer, real, etc.), e.g.
  . "print one integer as the answer".
- "array": The required output is a one-dimensional sequence (list) of values, e.g.
  an array of integers, a permutation, a sequence of answers for each query when
  printed as space-separated numbers or in multiple lines.
- "graph": The required output is a graph structure, such as a set of edges, tree
  description, adjacency list, or any structure where the output itself is
  naturally a graph.
- "matrix": The required output is a 2D grid or matrix.
- "bool": The required output is logically a boolean answer, e.g. "YES/NO", "True/
  False", "Alice/Bob", etc.
- "others": The required output is a complex or mixed structure that does not fit
  clearly into any of the above categories.
## 3. Output uniqueness
You must also decide whether the required output is unique for each valid input.
Define "is_output_unique" as:
- true if, for any fixed valid input, there is exactly one correct output that
  satisfies the problem statement.
- false if the statement allows multiple different outputs to be accepted as correct
  for the same input.
## JSON Output Specification
You must produce exactly one JSON object in a fenced JSON code block.
The JSON must have the following top-level keys:

```

```

- "scale_params": an object mapping parameter names to {{ "min": <int>, "max": <int>
  }}.
- "output_type": one of "string", "number", "array", "graph", "matrix", "bool", "
  others".
- "is_output_unique": a boolean.
## Example 1: <example 1>
## Example 2: <example 2>
## Example 3: <example 3>
## Final instruction
Now, read the provided problem statement and output the single JSON code block
  accordingly.
{problem}

```

917
918
919
920
921
922
923
924
925
926
927
928

A.2 Generate Test Case Prompt

930

Prompt 2: Prompt used for constructing generate_testcase function.

```

You are given as input a single algorithmic problem statement (like those from
  programming contests). Your job is to emit one Python code block that
  defines a test-case generator function for this problem.
The generator must produce exactly one valid test case per call, parameterized
  only by the numeric scale values provided via a JSON object.
## Input
You will be given:
1. A raw problem statement in natural language that fully specifies:
  - the input format,
  - the constraints,
  - and the meaning of each variable.
2. An example json_obj instance:
  - This is only an example to clarify field names and typical ranges.
  - Your code must work for any valid json_obj that matches the described schema.
## Required Python output (emit exactly one Python code block)
You must output a Python code block that defines one single function with the
  following signature:
```python
def generate_testcase(json_obj: dict) -> tuple[str, dict]:
 """
 Generate a test case based on the given json_obj.
 Parameters:
 - json_obj (dict): The input JSON object containing problem parameters.
 Returns:
 - tuple[str, dict]: A tuple containing:
 - The first element is a string representing the test case in input format.
 - The second element is a dictionary representing the same test case.
 """
 ...
 ...
Return value
- Your function must return both the string and the dictionary representation of the
 test case in a tuple. The first element of the tuple should be the string
 format, and the second element should be the dictionary format.
 - output_str:
 - A single string that is a valid input for the problem according to the
 Input section, representing exactly one logical test case.
 - If the problem statement defines a format with multiple test cases
 controlled by an integer T in the input, you must set T = 1.
 - Example: "1\n5\n1 2 3 4 5"
 - output_dict:
 - A Python dict that is a structured, formal description of the same test
 case.
 - If the problem statement contains multiple test cases, do not
 introduce T or any extra wrapper.
 - Example: {"n": 5, "list": [1, 2, 3, 4, 5]}
Constraints
- All sizes (counts, lengths, number of operations, etc.) must be determined only
 from json_obj.
- All other values (elements of arrays, weights, edges, indices, etc.) must be
 generated randomly within a reasonable range and strictly smaller than
 10000, while satisfying the problem's constraints at the same time.

```

931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983

```

984 - If the problem allows "no-solution" cases (e.g., the intended output is -1 when no
985 solution exists), you should **bias your random generation towards test cases
986 that admit at least one valid solution**, and explicitly construct values to
987 satisfy any hidden feasibility conditions, so that the correct solution is not
988 trivially always the "no-solution" output.
989 ## Problem statement
990 {problem}
991 ## Example json_obj
992 {example_json_obj}

```

## B Environment Cases

### B.1 Example Problem: 33\_C. Wonderful Randomized Sum

<b>Name</b>	33_C. Wonderful Randomized Sum
<b>Problem Description</b>	Learn, learn and learn again — Valera has to do this every day. He is studying at mathematical school, where math is the main discipline. The mathematics teacher loves her discipline very much and tries to cultivate this love in children. That's why she always gives her students large and difficult homework. Despite that Valera is one of the best students, he failed to manage with the new homework. That's why he asks for your help. He has the following task. A sequence of $n$ numbers is given. A prefix of a sequence is the part of the sequence (possibly empty), taken from the start of the sequence. A suffix of a sequence is the part of the sequence (possibly empty), taken from the end of the sequence. It is allowed to sequentially make two operations with the sequence. The first operation is to take some prefix of the sequence and multiply all numbers in this prefix by $-1$ . The second operation is to take some suffix and multiply all numbers in it by $-1$ . The chosen prefix and suffix may intersect. What is the maximum total sum of the sequence that can be obtained by applying the described operations?
<b>Generate Testcase</b>	<pre> def generate_testcase(json_obj, output_format="str"):     # Generate random test case for the problem     n = int(json_obj.get("n", 10))     numbers = random.sample(range(1, 100), n)      if output_format == "dict":         return {"n": n, "numbers": numbers}     else:         return f"{n}\n{' '.join(map(str, numbers))}" </pre>
<b>Output Requirement</b>	The first and the only line of the output should contain the answer to the problem.
<b>Difficulty Mapping</b>	{ "0": 0, "1": 1, "2": 2, "3": 2, "4": 3, "5": 4, "6": 5, "7": 6, "8": 8, "9": 11, "10": 14, "11": 18, "12": 23, "13": 30, "14": 39, "15": 51, "16": 67, "17": 87, "18": 112, "19": 146, "20": 190, "21": 247, "22": 318 }

### B.2 Example Problem: 1497\_D. Genius

<b>Name</b>	1497_D. Genius
-------------	----------------

<b>Problem Description</b>	Please note the non-standard memory limit. There are $n$ problems numbered with integers from 1 to $n$ . $i$ -th problem has the complexity $c_i = 2^i$ , tag $tag_i$ and score $s_i$ . After solving the problem $i$ it's allowed to solve problem $j$ if and only if $IQ <  c_i - c_j $ and $tag_i \neq tag_j$ . After solving it your IQ changes and becomes $IQ =  c_i - c_j $ and you gain $ s_i - s_j $ points. Any problem can be the first. You can solve problems in any order and as many times as you want. Initially your $IQ = 0$ . Find the maximum number of points that can be earned.
<b>Generate Testcase</b>	<pre>def generate_testcase(json_obj, output_format="str"):     # Generate a test case for the Genius problem     n = int(json_obj.get("n", 10))     tags = random.sample(range(1, n+1), n)     scores = random.sample(range(1, 100), n)      if output_format == "dict":         return {"n": n, "tags": tags, "scores": scores}     else:         return f"{n}\n{' '.join(map(str, tags))}\n{' '.join(map(str, scores))}"</pre>
<b>Output Requirement</b>	For each test case print a single integer — the maximum number of points that can be earned.
<b>Difficulty Mapping</b>	{ "0": 0, "1": 2, "2": 3, "3": 4, "4": 7, "5": 10, "6": 17, "7": 27, "8": 43, "9": 69, "10": 110, "11": 176, "12": 281, "13": 450, "14": 721, "15": 1153 }

### B.3 Example Problem: 1466\_B. Last Minute Enhancements

997

<b>Name</b>	1466_B. Last Minute Enhancements
<b>Problem Description</b>	Athenaeus has just finished creating his latest musical composition and will present it tomorrow to the people of Athens. Unfortunately, the melody is rather dull and highly likely won't be met with a warm reception. His song consists of $n$ notes, which we will treat as positive integers. The diversity of a song is the number of different notes it contains. As a patron of music, Euterpe watches over composers and guides them throughout the process of creating new melodies. She decided to help Athenaeus by changing his song to make it more diverse. Being a minor goddess, she cannot arbitrarily change the song. Instead, for each of the $n$ notes in the song, she can either leave it as it is or increase it by 1. Given the song as a sequence of integers describing the notes, find out the maximal, achievable diversity.
<b>Generate Testcase</b>	<pre>def generate_testcase(json_obj, output_format="str"):     # Generate test case for the song diversity problem     n = int(json_obj.get("n", 5))     notes = random.sample(range(1, 10), n)      if output_format == "dict":         return {"n": n, "notes": notes}     else:         return f"{n}\n{' '.join(map(str, notes))}"</pre>
<b>Output Requirement</b>	For each test case, you should output a single line containing precisely one integer, the maximal diversity of the song, i.e. the maximal possible number of different elements in the final sequence.

<b>Difficulty Mapping</b>	{ "0": 0, "1": 1, "2": 2, "3": 2, "4": 3, "5": 4, "6": 5, "7": 6, "8": 8, "9": 11, "10": 14, "11": 18, "12": 23, "13": 30, "14": 39, "15": 51, "16": 67, "17": 87, "18": 112, "19": 146, "20": 190, "21": 247, "22": 321, "23": 418 }
---------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

998

## C Training Detail

<b>Model Name</b>	Qwen3-4B-Base Qwen3-1.7B-Base
$K_{slope}$	10
$K_{zero}$	5
$K_{sat}$	5
<b>Learning Rate</b>	1e-6
<b>Learning Rate Warmup Steps</b>	20
<b>Batch Size</b>	64
<b>Max Prompt Length</b>	4096
<b>Max Response Length</b>	8192
<b>Entropy Coefficient</b>	0
<b>Number of Environments per Step</b>	64
<b>Training Prompt Batch Size</b>	64
<b>Mini Batch Size</b>	64
<b>Reward Estimator</b>	grpo
<b>KL Loss Coefficient</b>	0.0
<b>Clip Ratio Low</b>	0.2
<b>Clip Ratio High</b>	0.2
<b>Temperature (Training)</b>	1.0
<b>Top P (Training)</b>	1.0
<b>Top K (Training)</b>	-1
<b>Validation Temperature</b>	0.6
<b>Validation Top P</b>	0.95
<b>Number of Response per Prompt</b>	8