# Guided Self-Evolving LLMs with Minimal Human Supervision

**Wenhao Yu[1], Zhenwen Liang[1], Chengsong Huang[2], Kishan Panaganti[1],**
**Tianqing Fang[1], Haitao Mi[1], Dong Yu[1]**
[1]Tencent AI Lab in Seattle, [2]Washington University in St. Louis
wenhaowyu@global.tencent.com

## Abstract

AI self-evolution has long been envisioned as a path toward superintelligence, where models autonomously acquire, refine, and internalize knowledge from their own learning experiences. Yet in practice, unguided self-evolving systems often plateau quickly or even degrade as training progresses. These failures arise from issues such as concept drift, diversity collapse, and mis-evolution, as models reinforce their own biases and converge toward low-entropy behaviors. To enable models to self-evolve in a stable and controllable manner while minimizing reliance on human supervision, we introduce R-Few, a guided Self-Play Challenger–Solver framework that incorporates lightweight human oversight through in-context grounding and mixed training. At each iteration, the Challenger samples a small set of human-labeled examples to guide synthetic question generation, while the Solver jointly trains on human and synthetic examples under an online, difficulty-based curriculum. Across math and general reasoning benchmarks, R-Few achieves consistent and iterative improvements. For example, Qwen3-8B-Base improves by +3.0 points over R-Zero on math tasks and achieves performance on par with General-Reasoner, despite the latter being trained on 20 times more human data. Ablation studies confirm the complementary contributions of grounded challenger training and curriculum-based solver training, and further analysis shows that R-Few mitigates drift, yielding more stable and controllable co-evolutionary dynamics.
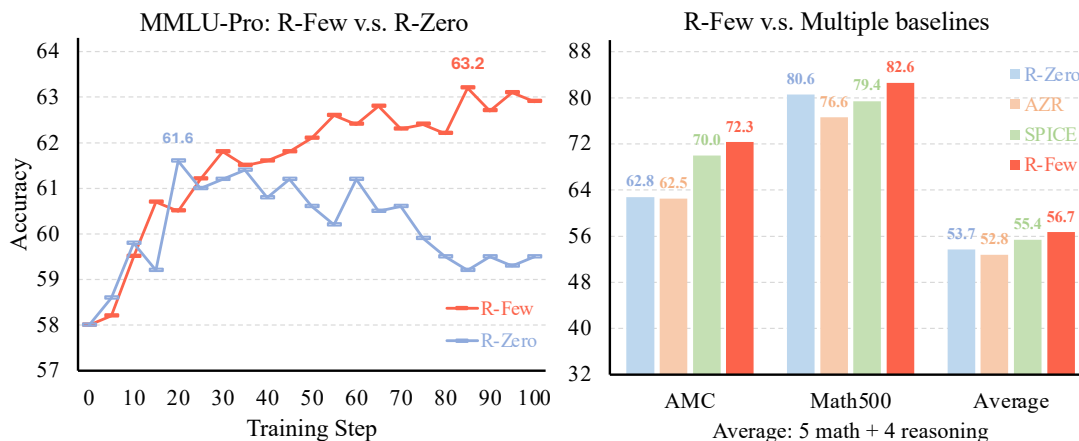


Figure 1: R-Few delays the performance plateau seen in R-Zero and achieves higher performance. After training, it outperforms baselines multiple benchmarks, showing more stable self-evolution.

## 1 Introduction

A long-standing vision in AI research is the development of systems that can self-evolve. It refers to a model that can autonomously acquire, refine, and internalize knowledge from its own generated experiences, thereby achieving continuous self-improvement without relying on large human-labeled datasets (Tao et al., 2024). To scale these capabilities without human supervision, self-play

offers a promising paradigm (Silver et al., 2017a; Sukhbaatar et al., 2017), where models improve by competing against themselves and generating automatic feedback through competition.

However, when applied to language tasks, self-play remains surprisingly brittle. Despite early gains, methods such as Absolute Zero and R-Zero often plateau quickly, and their performance may even deteriorate as training progresses. (Zhao et al., 2025a; Huang et al., 2025; Chen et al., 2025; Kuba et al., 2025; He et al., 2025). Two key challenges underlie this phenomenon: (i) *Concept Drift*: Without external intervention or guidance, models tend to reinforce their own knowledge bias. Over repeated self-play iterations, feedback loops amplify spurious correlations or biased reasoning patterns, drifting away from factual correctness and semantic validity (Liu et al., 2025a). (ii) *Diversity Collapse*: Since the underlying knowledge of the model is fixed, its self-generated challenges tend to converge toward familiar and low-entropy regions of the task space. As the diversity of generated tasks diminishes, exploration and reasoning novelty stagnate (Liang et al., 2025; Morris et al., 2025).

Together, these problems make it difficult for language models to evolve in a stable and controllable manner. To overcome this challenge, we need a framework that minimizes reliance on human annotations while its evolution should remain guided toward human-intended goals. In this work, we propose R-FEW, a self-evolving LLM framework that integrates limited human supervision into the self-play process. Unlike prior approaches that require large-scale labeled corpora or extensive human feedback, R-FEW leverages only a small pool of high-quality "anchor" data – specifically, 1% to 5% data sampled from WebInstruct (Ma et al., 2025), a open-source large-scale dataset mined from the web. During training, the Challenger dynamically samples zero to five anchor examples as in-context example to guide synthetic question generation. When zero examples are sampled, R-FEW preserves the ability for open-ended exploration; when few-shot examples are included, they inject semantic anchors that regularize the Challenger's generation. To further enhance the learning dynamics, we introduce an online curriculum mechanism for the Solver. Instead of indiscriminately training on all generated data, R-FEW first rolls out multiple solving trajectories, ranks questions by the Solver's uncertainty, and selects mid-uncertain samples for the next iteration training. This ranking-based curriculum ensures that the Solver focuses on the zone of proximal development. The same selection principle is applied to both synthetic and human anchor data, so it can merge human-curated and self-generated question into a unified training stream for Solver.

Our experiments demonstrate that R-FEW consistently and iteratively improving the reasoning abilities. For example, Qwen3-8B-Base model's average score on math benchmarks increased by a significant +3.0 points than R-Zero, achieving on par performance with General-Reasoner (Ma et al., 2025), the model trained on the full WebInstruct dataset (232k). In addition, we provide an in-depth analysis that validates our framework's components, and characterizes the co-evolutionary dynamics to identify both strengths and limitations, offering insights for future research.

## 2 Preliminaries: Self-Play For Data-Free Training

Inspired by the success of self-play mechanisms in games, exemplified by AlphaZero (Silver et al., 2017a;b), recent research have begun extending this paradigm to large language models (LLMs), aiming to create self-evolving systems that can continually refine their reasoning abilities without curated human data (Chen et al., 2024). As base model capabilities have risen and reinforcement learning has advanced, a wave of "self-evolving + RL" training schemes has emerged (Zhao et al., 2025a; Huang et al., 2025; Chen et al., 2025; Liu et al., 2025a). These approaches show that models can autonomously construct their own curricula – posing challenges, solving them, and improving through iteration – making self-improving intelligence an increasingly attainable reality.

### 2.1 Unified View of Self-Play Objectives and Reward Design

Let $Q_\theta$ denote the **Challenger** (aka. Proposer, Questioner) and $S_\phi$ the **Solver** (aka. Reasoner), both initialized from the same pretrained LLM $M_0$. At each iteration $t$, the Challenger samples a question or task $q_t \sim Q_\theta(\cdot \mid \mathcal{H}_t)$, where $\mathcal{H}_t$ represents the historical context (e.g., previous questions and responses, document corpus), and the Solver attempts to generate an answer $a_t \sim S_\phi(\cdot \mid q_t)$. A verifier or self-consistency mechanism then evaluates outputs and produces a scalar reward $r_t = \mathcal{R}(q_t, a_t; S_\phi)$, which guides the policy updates of both agents via reinforcement learning,

typically using Group Relative Policy Optimization (GRPO) (Shao et al., 2024):

$$\theta \leftarrow \theta + \eta_\theta \nabla_\theta \mathbb{E}_{q_t \sim Q_\theta}[\mathcal{R}_{\text{chal}}(q_t)], \quad (1)$$

$$\phi \leftarrow \phi + \eta_\phi \nabla_\phi \mathbb{E}_{a_t \sim S_\phi(q_t)}[\mathcal{R}_{\text{sol}}(q_t, a_t)]. \quad (2)$$

**Unified objective.** The overall bi-level optimization can be expressed as

$$\max_\theta \mathbb{E}_{q \sim Q_\theta}[\mathcal{R}_{\text{chal}}(q; \phi)] \quad \text{s.t.} \quad \phi \leftarrow \arg\max_\phi \mathbb{E}_{q \sim Q_\theta, a \sim S_\phi(q)}[\mathcal{R}_{\text{sol}}(q, a)]. \quad (3)$$

This formulation unifies language self-play methods under a single reinforcement learning objective, where the Challenger adaptively generates a curriculum of problems and the Solver improves its reasoning skills in response.

**Reward design across frameworks.** Different self-play frameworks instantiate distinct choices of the Challenger and Solver rewards ($\mathcal{R}_{\text{chal}}, \mathcal{R}_{\text{sol}}$), each encoding a specific learning bias:

- **Absolute Zero** (Zhao et al., 2025a): Operates in a fully zero-data regime, using a linear difficulty shaping function to maintain stability:

$$\mathcal{R}_{\text{chal}}^{\text{AbsZero}}(q) = (1 - \widehat{p}_{\text{succ}}(q)) \cdot \mathbb{I}(\widehat{p}_{\text{succ}}(q) \neq 0), \qquad \mathcal{R}_{\text{sol}}^{\text{AbsZero}}(q, a) = \mathbb{I}\{a = y(q)\},,$$

  where $f$ peaks near $\widehat{p}_{\text{succ}} = 0$

- **R-Zero** (Huang et al., 2025): Employs a purely uncertainty-driven curriculum with verifiable pseudo-labels:

$$\mathcal{R}_{\text{chal}}^{\text{R-Zero}}(q) = 1 - 2\left|\widehat{p}_{\text{succ}}(q) - \tfrac{1}{2}\right| - \lambda_{\text{rep}}\text{RepPenalty}(q), \qquad \mathcal{R}_{\text{sol}}^{\text{R-Zero}}(q, a) = \mathbb{I}\{a = \tilde{y}(q)\},$$

  where $\tilde{y}(q)$ is the majority-vote pseudo-answer and RepPenalty discourages duplicate questions.

- **Self-Questioning LMs (SQLM)** (Chen et al., 2025): The Challenger optimizes for both *answerability* and *informativeness* by targeting a difficulty sweet spot based on the success rate statistics:

$$\mathcal{R}_{\text{chal}}^{\text{SQLM}}(q) = (1 - \widehat{p}_{\text{succ}}(q)) \cdot \mathbb{I}(0 < p_{\text{succ}}(q) < \widehat{p}_{\text{succ}}(q)), \qquad \mathcal{R}_{\text{sol}}^{\text{SQLM}}(q, a) = \mathbb{I}\{a = \tilde{y}(q)\},$$

  where $p_{\text{succ}}(q)$ denotes the pass rate, and $\tilde{y}(q)$ represents the majority-vote pseudo-answer.

- **SPICE** (Liu et al., 2025a): Uses a *corpus-guided* framework where the Challenger is rewarded for generating questions with a specific response variance (indicating appropriate difficulty), while the Solver is supervised by the source text:

$$\mathcal{R}_{\text{chal}}^{\text{SPICE}}(q) = \mathbb{I}(q \text{ is valid}) \cdot \exp\left(-\frac{(\text{Var}(\{\hat{y}_1, \ldots, \hat{y}_K\}) - 0.25)^2}{2 \cdot 0.01}\right), \quad \mathcal{R}_{\text{sol}}^{\text{SPICE}}(q, a) = \mathbb{I}\{a = y(q)\},$$

  where $y(q)$ is the ground-truth answer extracted directly from the reference document, and $\{\hat{y}_k\}$ represents the set of sampled responses from peer models.

## 3 Method

### 3.1 Overview

Building upon the data-free self-play foundation established, we introduce R-FEW, a self-evolving framework that integrates *limited human supervision* to stabilize and accelerate reasoning evolution. R-FEW presents two key innovations: (1) a **few-shot grounded Challenger** that anchors synthetic task generation to small human "anchor" data, and (2) an **online curriculum Solver** that adaptively selects mid-uncertain samples from both synthetic and human sources.
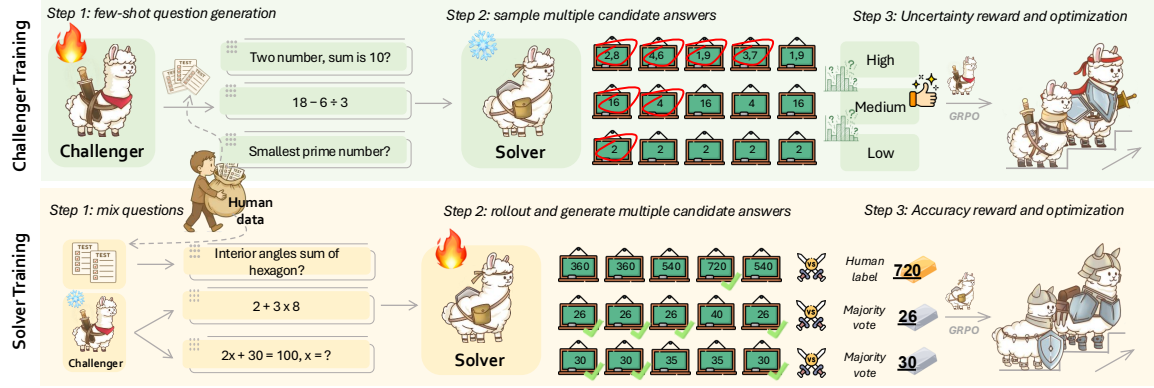
Figure 2: **The math examples shown in the figure are not real data**, but are included for demonstration to aid understanding. The figure provides an overview of our R-FEW framework. The Challenger is incentivized to generate moderately ("medium") uncertain questions that lie at the edge of the Solver's current abilities; the Solver is rewarded for solving increasingly challenging tasks – sourced from both humans and the Challenger – via curriculum-based selection.

## 3.2 Few-Shot Grounded Challenger

The Challenger $Q_\theta$ in R-FEW generates synthetic questions conditioned on a small pool of high-quality human anchor examples $\mathcal{D}_H = \{(x_i, y_i)\}_{i=1}^{N_H}$. At each rollout step, we randomly sample $k \in \{0, 1, \ldots, 5\}$ demonstrations as in-context examples:

$$\mathcal{C}_t = \text{Sample}_k(\mathcal{D}_H), \qquad q_t \sim Q_\theta(\cdot \mid \mathcal{C}_t, \mathcal{H}_t). \tag{4}$$

When $k = 0$, R-FEW degenerates to the data-free setting of R-ZERO, preserving open-ended self-evolution. For $k > 0$, the Challenger is softly guided toward human-consistent reasoning trajectories while retaining generative diversity.

$$\mathcal{R}_{\text{chal}}^{\text{R-FEW}}(q_t) = \underbrace{1 - 2\left|\widehat{p}_{\text{succ}}(q_t) - \tfrac{1}{2}\right|}_{\text{difficulty shaping}} - \lambda_{\text{rep}} \text{RepPenalty}(q_t). \tag{5}$$

Here, $\text{Align}(q_t, \mathcal{D}_H)$ measures the semantic or structural proximity between the generated question $q_t$ and the human anchor distribution via cosine similarity in embedding space. This term encourages the Challenger to explore *in the vicinity* of real human data, preventing semantic drift while maintaining autonomous difficulty progression.

The Challenger parameters are updated via GRPO as:

$$\theta \leftarrow \theta + \eta_\theta \nabla_\theta \mathbb{E}_{q_t \sim Q_\theta}[\mathcal{R}_{\text{chal}}^{\text{R-FEW}}(q_t)]. \tag{6}$$

Besides, there is another trick we used during challenger training. Because the base model sometimes struggles to follow instructions when the prompt gets longer, we added a quick warm-up stage via supervised fine-tuning (SFT) to help it reliably adhere to the required format.

## 3.3 Online Curriculum Solver

To ensure efficient learning and avoid collapse, the Solver $S_\phi$ employs an **online curriculum mechanism** that dynamically ranks problems by difficulty and filters the training set to the "zone of proximal learning". At iteration $t$, the Challenger produces a batch of $K$ questions $\{q_t^{(i)}\}_{i=1}^K$. For

---

**Algorithm 1: R-FEW**

---

**Require:** Pretrained base LLM $M_0$; human anchor data $\mathcal{D}_H$; batch size $B$; group size $G$;
iterations $T$; anchor sampling range $k \in [0, 5]$; curriculum quantiles $[\tau_{\text{low}}, \tau_{\text{high}}]$.
**Return:** Trained Challenger $Q_\theta$ and Solver $S_\phi$

---

1　**for** $t \leftarrow 1$ **to** $T$ **do**
　　// Challenger Role: Few-Shot Grounded Generation
2　　**for** $b \leftarrow 1$ **to** $B$ **do**
3　　　　Sample $k$ anchor examples: $\mathcal{C}_t \leftarrow \text{Sample}_k(\mathcal{D}_H)$;
4　　　　Generate $G$ question candidates via few-shot learning: $\{q_i\}_{i=1}^G \sim Q_\theta(\cdot \mid \mathcal{C}_t, \mathcal{H}_t)$;
5　　　　Evaluate Solver's success probabilities: $\hat{p}_{\text{succ}}(q_i) \leftarrow \frac{1}{M}\sum_{m=1}^M J(q_i, a_i^{(m)})$;
6　　　　**if** $q_i$ is valid **then**
7　　　　　　Compute Challenger reward:

$$r_C(q_i) \leftarrow \big(1 - 2|\hat{p}_{\text{succ}}(q_i) - 0.5|\big) - \lambda_{\text{rep}}\,\text{RepPenalty}(q_t).$$

8　　　　**else**
9　　　　　　$r_C(q_i) \leftarrow -\rho_{\text{inv}} \cdot \mathbf{1}[\text{invalid}]$;
10　　Update Challenger:
$$\theta \leftarrow \theta + \eta_\theta \nabla_\theta \mathbb{E}_{q_i \sim Q_\theta}[r_C(q_i)] \quad \text{via GRPO.}$$

　　// Solver Role: Online Curriculum Learning
11　　Aggregate question–answer pairs $\{(q_i, a_i)\}$ from Challenger and human anchors $\mathcal{D}_H$;
12　　Compute success rates $\hat{p}_{\text{succ}}(q)$ for all $q \in \{q_i\} \cup \mathcal{D}_H$;
13　　Select mid-difficulty subset: $\mathcal{D}_{\text{cur}} \leftarrow \{(q, a) : \tau_{\text{low}} \leq \hat{p}_{\text{succ}}(q) \leq \tau_{\text{high}}\}$;
14　　Set $\mathcal{D}_{\text{mix}} \leftarrow \mathcal{D}_{\text{cur}}$;
15　　**for** $(q, a) \in \mathcal{D}_{\text{mix}}$ **do**
16　　　　Compute Solver reward:

$$r_S(q, a) \leftarrow w_{\text{cur}}(q)\mathbf{1}[a = \tilde{y}(q)] + \lambda_{\text{hum}} w_{\text{hum}}(q)\mathbf{1}[(q, a) \in \mathcal{D}_H]$$

17　　Update Solver:
$$\phi \leftarrow \phi + \eta_\phi \nabla_\phi \mathbb{E}_{(q,a) \sim \mathcal{D}_{\text{mix}}}[r_S(q, a)] \quad \text{via GRPO.}$$

18　**return** Trained Challenger $Q_\theta$ and Solver $S_\phi$

---

each $q_t^{(i)}$, the Solver performs $M$ rollouts and records the success rate $\hat{p}_{\text{succ}}(q_t^{(i)})$:

$$\hat{p}_{\text{succ}}(q_t^{(i)}) = \frac{1}{M}\sum_{m=1}^M J\left(q_t^{(i)}, a_t^{(i,m)}\right), \qquad a_t^{(i,m)} \sim S_\phi(\cdot \mid q_t^{(i)}). \tag{7}$$

We then sort all questions (both synthetic and human) by success rate and select mid-range items according to the quantile interval $[\tau_{\text{low}}, \tau_{\text{high}}]$, empirically we set it as $[0.3, 0.7]$:

$$\mathcal{D}_{\text{cur}} = \big\{(q, a) : \tau_{\text{low}} \leq \hat{p}_{\text{succ}}(q) \leq \tau_{\text{high}}\big\}. \tag{8}$$

This ensures the Solver focuses on problems that are challenging yet solvable, consistent with optimal learning theory (Shi et al., 2025; Bae et al., 2025). The Solver reward incorporate both synthetic and human data, weighted by curriculum difficulty:

$$\mathcal{R}_{\text{sol}}^{\text{R-FEW}}(q, a) = w_{\text{cur}}(q) \cdot \mathbf{1}\{a = \tilde{y}(q)\} + \lambda_{\text{hum}} w_{\text{hum}}(q) \cdot \mathbf{1}\{(q, a) \in \mathcal{D}_H\}, \tag{9}$$

where $w_{\text{cur}}(q)$ is a normalized curriculum weight, and $w_{\text{hum}}(q)$ upweights the scarce human anchors to prevent forgetting. In practice, we set $\lambda$ as 2.0. Solver parameters are optimized as:

$$\phi \leftarrow \phi + \eta_\phi \nabla_\phi \, \mathbb{E}_{(q,a) \sim \mathcal{D}_{\text{cur}}}[\mathcal{R}_{\text{sol}}^{\text{R-FEW}}(q,a)]. \tag{10}$$

### 3.4 Iterative Co-Evolution

R-FEW maintains the same self-evolving loop as R-ZERO, but now with a grounded and curriculum-aware data pipeline: Challenger generates synthetic tasks with few-shot demonstrations from human anchors. Solver attempts to solve tasks and estimates success rates. An online curriculum filter selects mid-difficulty tasks from both synthetic and human pools.

## 4 Experiments

### 4.1 Experiments Setting

#### 4.1.1 Baseline Methods

We employ Qwen3-4B-Base (Yang et al., 2025b) and Qwen3-8B-Base as our backbone models. For comparison, we evaluate against several baselines:
(1) Base Model — the pretrained checkpoint without any post-training, and the model checkpoint can be downloaded at `https://huggingface.co/Qwen/Qwen3-4B-Base` (`Qwen3-8B-Base`);
(2) R-Zero (Huang et al., 2025) — the first ungrounded self-play approach where the model generates its own questions and answers, enabling self-improvement through autonomous reasoning;
(3) Absolute Zero (Zhao et al., 2025a) — a self-play framework constrained to code generation tasks with Python execution as verification, exemplifying domain-specific grounded self-play;
(4) SPICE (Liu et al., 2025a) — a grounded self-play framework that retrieves contexts from a large document corpus to create diverse reasoning tasks with document-grounded answers.

We reproduce baselines (1)-(2) using publicly available implementations under identical training infrastructure, while results for (3)–(4) are taken from the corresponding reports in Liu et al. (2025a).

We evaluate two variants of our method—R-FEW (1%) and R-FEW (5%)—corresponding to using 1% or 5% of examples randomly sampled from the WebInstruct dataset (232k) (Ma et al., 2025).

#### 4.1.2 Evaluation Benchmark

We assess our framework on a comprehensive suite of benchmarks under two main categories.

**Mathematical Reasoning.** We use five benchmarks: AMC, Minerva (Lewkowycz et al., 2022), MATH-500 (Hendrycks et al., 2021), GSM8K (Cobbe et al., 2021) and Olympiad-Bench (He et al., 2024). For evaluation, we follow Zhao et al. (2025c); Ma et al. (2025) to employ GPT-4o as a programmatic judge to semantically verify the correctness of the final answer against the ground truth.

**General Domain Reasoning.** To test for the generalization of reasoning ability, we evaluate on the following benchmarks: MMLU-Pro (Wang et al., 2024), SuperGPQA (Du et al., 2025), GPQA-Diamond (Rein et al.), and BBEH (shoaa kazemi et al., 2025). These benchmarks collectively cover a wide range of disciplines (e.g., physics, biology, business, economic, law) with challenging tasks, offering a stricter evaluation of complex reasoning abilities. For evaluation, we follow Ma et al. (2025) to report Exact Match (EM) accuracy obtained via greedy decoding.

### 4.2 Experimental Results

#### 4.2.1 Comparison with Baseline Methods

Table 1 demonstrates that R-FEW enables strong and efficient self-evolution across both mathematical and general reasoning tasks. For Qwen3-4B-Base, the Base Model achieves an average score of 41.9, and unguided self-play via R-Zero improves performance to 48.2. However, these gains remain limited, illustrating the difficulty of stable self-evolution without human grounding. In contrast,

Table 1: Comprehensive results on reasoning benchmarks. We compare our R-Few against several self-evolving baselines, including R-Zero, Absolute Zero, SPICE. R-Few consistently improves performance over baselines, and approaching General-Reasoner (trained with 232k WebInstruct data) while using substantially less human supervision (1% and 5%).

| Models ↓ | Mathmetical Reasoning | | | | | General Reasoning | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AMC | Minerva | MATH 500 | GSM8K | Olympiad | MMLU Pro | Super GPQA | GPQA-D | BBEH | Avg. |
| *Qwen3-4B-Base* | | | | | | | | | | |
| Base Model | 47.5 | 42.3 | 68.2 | 72.6 | 34.8 | 51.6 | 25.4 | 26.3 | 8.1 | 41.9 |
| + R-Zero | 48.2 | 51.2 | 74.8 | 90.6 | 40.6 | 54.2 | 27.8 | 36.4 | 10.4 | 48.2 |
| + Absolute Zero | 50.0 | 41.9 | 76.2 | 89.3 | 41.5 | 52.6 | 27.1 | 35.3 | 8.3 | 46.9 |
| + SPICE | 57.5 | 51.9 | 78.0 | 92.7 | 42.7 | 58.1 | 30.2 | 39.4 | 12.3 | 51.4 |
| + R-Few (1%) | 52.7 | 52.1 | 77.8 | 92.3 | 42.4 | 55.9 | 29.4 | 35.4 | 11.2 | 49.9 |
| + R-Few (5%) | 52.4 | 53.2 | 78.0 | 92.6 | 42.8 | 56.2 | 29.4 | 39.9 | 11.8 | 50.7 |
| General-Reasoner | 60.0 | 57.7 | 80.6 | 92.2 | 47.7 | 62.8 | 32.5 | 42.9 | 12.2 | 54.3 |
| *Qwen3-8B-Base* | | | | | | | | | | |
| Base Model | 61.5 | 49.3 | 74.4 | 90.9 | 40.4 | 58.0 | 30.4 | 33.3 | 10.5 | 49.9 |
| + R-Zero | 62.8 | 58.8 | 80.6 | 92.4 | 43.4 | 61.6 | 31.8 | 40.5 | 11.3 | 53.7 |
| + Absolute Zero | 62.5 | 52.9 | 76.6 | 92.0 | 47.8 | 62.5 | 33.5 | 36.8 | 10.8 | 52.8 |
| + SPICE | <u>70.0</u> | 59.2 | 79.4 | 92.7 | 42.5 | 65.0 | 35.7 | 39.4 | **14.9** | 55.4 |
| + R-Few (1%) | 69.3 | 59.6 | 81.6 | **94.0** | 44.0 | 62.8 | 32.7 | 40.4 | 11.8 | 55.1 |
| + R-Few (5%) | **72.3** | <u>60.3</u> | <u>82.6</u> | 93.5 | **46.4** | <u>63.2</u> | <u>33.5</u> | **46.5** | <u>12.3</u> | **56.7** |
| General-Reasoner | 64.8 | **62.6** | **83.4** | 92.7 | <u>46.3</u> | **65.1** | **35.3** | <u>42.9</u> | 10.8 | <u>56.0</u> |

R-FEW produces significantly larger improvements: with only 1% human data, performance rises to 49.9, and with 5% data it reaches 50.7, narrowing the gap to General-Reasoner (54.3), which relies on roughly 20× more human-labeled data. This demonstrates that minimal human grounding can unlock meaningful self-evolution while retaining data efficiency.

A notable trend emerges when comparing model scales. The 8B variants exhibit stronger self-evolving capability than their 4B counterparts. While R-Zero provides modest improvements for both, R-FEW on Qwen3-8B-Base reaches 55.1 with 1% data and 56.7 with 5%, surpassing General-Reasoner (56.0). This indicates that larger models not only benefit more from guided self-evolution but can also nearly replicate or exceed the performance of heavily supervised pipelines with minimal human input. These observations suggest that as model capacity increases, the model becomes more capable of interpreting human grounding signals, generating higher-quality synthetic data, and improving its reasoning abilities through iterative self-play.

### 4.2.2 Ablation Studies

To isolate the contribution of each key component based on our R-FEW (5%) experiments, we conduct a comprehensive ablation study on the `Qwen3-8B-Base` model. We evaluate the importance of three critical modules by disabling each one individually and measuring its impact on math and general reasoning tasks. The results are summarized in Table 2.

Among the ablations, disabling challenger training causes the largest degradation, reducing Math and General averages by 1.9 and 1.0 points, respectively. Removing challenger warm-up and curriculum learning yields similar trends, indicating both components are critical for stronger performance. Math performance is more sensitive than general domain performance. This sensitivity is consistent with the trend observed in Table 1, where math benchmarks exhibit larger performance variance across different methods.

Table 2: Ablation study on the Qwen3-8B-Base model. Results highlight the impact of challenger training, challenger warm-up, and curriculum learning on downstream reasoning tasks.

| Method | Math | General |
|---|---|---|
| R-FEW (from Table 1) | 71.0 | 38.9 |
| *Ablations* | | |
| ⊢ w/o Challenger Training | 68.0 | 37.4 |
| ⊢ w/o Challenger Warm-up | 69.1 | 37.9 |
| ⊢ w/o Curriculum Learning | 69.1 | 38.1 |

### 4.2.3 Domain Corelation with Human Data

Figure 3 presents the cross-domain performance matrix, showing how human-generated data from each training category influences performance across MMLU-Pro subcategories. We sample domain-specific data from WebInstruct (already labeled) to train the self-evolving models, and then evaluate them on MMLU-Pro. A clear pattern emerges: data from each domain tends to yield the greatest improvements within its own category. Notably, math stands out as the most useful category, suggesting that mathematical reasoning contributes broadly to complex reasoning capabilities. The cross-domain relationships also highlight two strongly connected pairs: math–physics, reflecting shared quantitative structure, and business–economics, likely because of their overlapping analytical and decision-making frameworks.
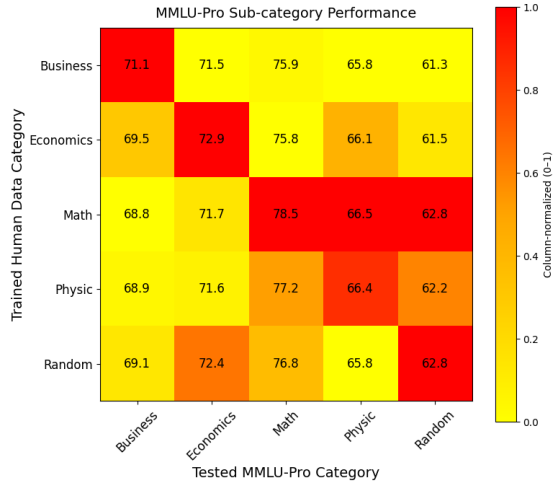


Figure 3: Impact of domain-sampled human data on performance across MMLU-Pro categories.

Overall, the figure indicates that while heterogeneous human data provides general improvements, domain-specific sampling is particularly effective at strengthening the domain it represents. This observation is valuable because it makes the self-evolving system more controllable: by selecting or emphasizing particular data domains, developers can steer how the model grows, prevent unintended concept drift, and maintain clearer alignment between training inputs and emergent skills, rather than drifting unpredictably into unrelated areas.

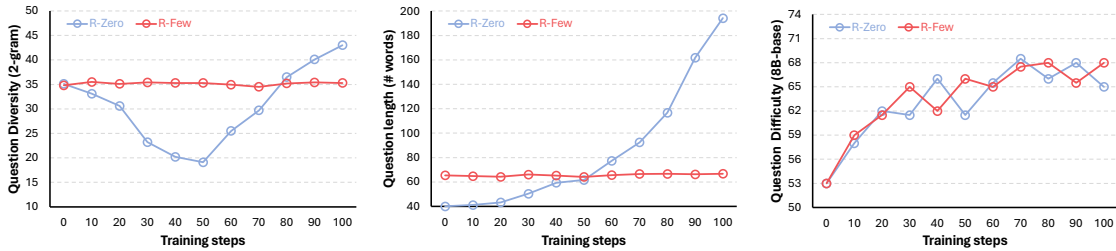### 4.2.4 Improved Stability and Reduced Reward Hacking



Figure 4: Training curves of synthetic question diversity (measured by 2-gram lexical diversity), length (measured by word count), and difficulty (evaluated by Qwen3-8B-Base, with ground-truth labeled by Gemini-2.5-Pro) over training. R-Zero collapses in diversity and exhibits length inflation via verbosity, whereas R-FEW maintains stable length and diversity during self-evolution.

As discussed in the introduction, unguided self-evolving systems tend to plateau quickly and exhibit instability, often leading to diversity collapse and reward hacking. Our empirical results clearly reflect these phenomena. In the left panel, R-Zero shows a sharp decline in question diversity as training progresses, dropping from 35 to below 20 in the first 50 training steps. Although the 2-gram diversity score increases after 50 steps, this increase coincides with a dramatic explosion in question length, as shown in the middle panel, indicating that is largely a length-induced artifact rather than genuine semantic diversification. This is because longer outputs naturally introduce more unique n-gram combinations, inflating the diversity metric without reflecting meaningful exploration. A similar pattern emerges in question length itself. R-Zero progressively inflates question length, eventually producing excessively long questions. This suggests reward hacking, where the model exploits superficial features, such as verbosity, to increase perceived diversity or difficulty signals. In contrast, R-FEW maintains a stable diversity throughout training while keeping question length consistent, suggesting that lightweight human grounding prevents diversity collapse.

Finally, the right panel shows that both methods increase question difficulty over time, but this similarity masks important qualitative differences. To measure difficulty, we relabel questions

8

generated at different training stages using Gemini-2.5-Pro, then evaluate them by having the same model (Qwen3-8B-Base) attempt to answer them. Difficulty is defined as the percentage of questions the model answers incorrectly (the inverse of accuracy). We observe that both curves exhibit an overall upward trend, indicating increasing difficulty; however, the source of this increase differs substantially. R-Zero produces questions that are harder to solve primarily because they are longer, rather than because they require deeper reasoning. This reflects another form of reward hacking, where verbosity inflates perceived difficulty without meaningful improvement in reasoning complexity. In contrast, R-FEW maintains consistent question length, allowing the model to self-evolve toward genuinely more challenging questions without sacrificing diversity or structure.

## 5 Related Work

### 5.1 Self-Evolving LLMs and Self-Play

Early demonstrations of self-play–based learning came from AlphaZero, which learns without external supervision by playing against itself and progressively strengthening its policy through these matches (Silver et al., 2017a;b; Sukhbaatar et al., 2017). Recent research extended this principle to language leads naturally to self-evolving LLMs that refine their own reasoning capabilities over time (Tao et al., 2024). Several systems instantiate language self-play at scale. Self-play fine-tuning improves relatively weak LLMs by repeatedly generating and solving tasks without large labeled corpora (Chen et al., 2024). This approach has been particularly fruitful in verifiable domains like code generation, where a "Coder" agent's program is verified by a "Tester" agent's unit tests (Lin et al., 2025; Wang et al., 2025; Pourcel et al., 2025; Zhao et al., 2025a). Fully data-free methods such as R-Zero push this further: the model proposes questions, produces candidate solutions, and learns purely from internal verification signals (Huang et al., 2025; Chen et al., 2025; Kuba et al., 2025). Other approaches introduce limited seeds or structure, where the model starts from a small set of examples and learns to synthesize increasingly challenging problems around them (Fang et al., 2025; He et al., 2025; Liu et al., 2025a). R-FEW builds on this line of work but targets a different point on the supervision spectrum. This guided self-play design aims to preserve the open-ended exploration benefits of self-evolving LLMs while explicitly mitigating concept drift and diversity collapse.

### 5.2 Reinforcement Learning for LLM Reasoning

Reinforcement learning (RL) has become a central mechanism for improving LLM reasoning beyond supervised fine-tuning and instruction tuning, most prominently through reinforcement learning from human feedback (RLHF) and its variants (Kaufmann et al., 2024). Early RLHF pipelines such as InstructGPT (Ouyang et al., 2022) demonstrated that sequence-level preference optimization can substantially improve helpfulness, safety, and, in many cases, reasoning quality compared to pure supervised fine-tuning. Building on these foundations, recent reasoning-centric systems such as DeepSeek-R1 (DeepSeek-AI et al., 2025; Shao et al., 2024) show that RL-style optimization with carefully designed reasoning rewards can significantly enhance multi-step problem solving. Similar R1-style pipelines in both commercial and open-source models further corroborate that tailoring RL objectives to long-form chain-of-thought and verifiable tasks yields substantial gains in reasoning ability (Yang et al., 2025a; Team et al., 2025; Liu et al., 2025b; Yu et al., 2025).

Within this landscape, RL with verifiable rewards (RLVR) uses externally defined verifiers or environments to provide scalar feedback (Ma et al., 2025; Su et al., 2025; Yu et al., 2025). For instance, Search-R1 teaches LLMs to reason while leveraging search engines (Jin et al., 2025), and semantically-aware R1 variants extend RLVR to open-ended free-form generation (Li et al., 2025b). These works illustrate how verifiable environments – logical constraints, executable code, structured prediction targets, or retrieval-enhanced setups – can provide dense, automatically computable rewards beyond human labels. Another trend in recent research is label-free reinforcement learning, which aims to improve LLM reasoning without human-annotated data or explicit external verifiers. Many such methods use self-generated outputs as a reward signal. This includes leveraging sequence-level confidence (Li et al., 2025a; Prabhudesai et al., 2025), the consistency of answers derived from varied reasoning paths (Zhang et al., 2025; Zuo et al., 2025), or minimizing the output entropy (Agarwal et al., 2025; Cheng et al., 2025). These signals are often used within self-training loops where models fine-tune on their own most plausible solutions (Shafayat et al., 2025; Zhao et al., 2025b).

## 6 Conclusion and Future Work

In this paper, we introduced R-FEW, a guided self-evolving framework that enables LLMs to improve autonomously with minimal human supervision. By combining few-shot–grounded question generation with an online curriculum for solver training, R-FEW produces stable, scalable co-evolution and achieves strong gains across mathematical and general reasoning tasks. Despite using only 1–5% human data, it approaches or surpasses systems trained with far larger labeled datasets, demonstrating the power of lightweight anchoring in preventing drift and enhancing learning dynamics. Future work includes improving efficiency, exploring richer verification methods, and extending self-evolution to open-ended domains lacking objective correctness signals.

## References

Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han, and Hao Peng. The unreasonable effectiveness of entropy minimization in llm reasoning. *ArXiv preprint*, abs/2505.15134, 2025.

Sanghwan Bae, Jiwoo Hong, Min Young Lee, Hanbyul Kim, JeongYeon Nam, et al. Online difficulty filtering for reasoning oriented reinforcement learning. *ArXiv preprint*, abs/2504.03380, 2025.

Lili Chen, Mihir Prabhudesai, Katerina Fragkiadaki, Hao Liu, and Deepak Pathak. Self-questioning language models. *arXiv preprint arXiv:2508.03682*, 2025.

Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. In *International Conference on Machine Learning*, pp. 6621–6642. PMLR, 2024.

Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai, Wayne Xin Zhao, et al. Reasoning with exploration: An entropy perspective. *ArXiv preprint*, abs/2506.14758, 2025.

Karl Cobbe, Vineet Kosaraju, Mo Bavarian, Mark Chen, Heewoo Jun, et al. Training verifiers to solve math word problems. *ArXiv preprint*, abs/2110.14168, 2021.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Jun-Mei Song, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *ArXiv preprint*, abs/2501.12948, 2025.

Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, et al. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. *ArXiv preprint*, abs/2502.14739, 2025.

Wenkai Fang, Shunyu Liu, Yang Zhou, Kongcheng Zhang, Tongya Zheng, et al. Serl: Self-play reinforcement learning for large language models with limited data. *ArXiv preprint*, abs/2505.20347, 2025.

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Annual Meeting of the Association for Computational Linguistics*, 2024.

Yicheng He, Chengsong Huang, Zongxia Li, Jiaxin Huang, and Yonghui Yang. Visplay: Self-evolving vision-language models from images. *arXiv preprint arXiv:2511.15661*, 2025.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, et al. Measuring mathematical problem solving with the math dataset. *ArXiv preprint*, abs/2103.03874, 2021.

Chengsong Huang, Wenhao Yu, Xiaoyang Wang, Hongming Zhang, Zongxia Li, Ruosen Li, Jiaxin Huang, Haitao Mi, and Dong Yu. R-zero: Self-evolving reasoning llm from zero data. *arXiv preprint arXiv:2508.05004*, 2025.

Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *ArXiv preprint*, abs/2503.09516, 2025.

Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. 2024.

Jakub Grudzien Kuba, Mengting Gu, Qi Ma, Yuandong Tian, and Vijai Mohan. Language self-play for data-free training. *arXiv preprint arXiv:2509.07414*, 2025.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35: 3843–3857, 2022.

Pengyi Li, Matvey Skripkin, Alexander Zubrey, Andrey Kuznetsov, and Ivan V. Oseledets. Confidence is all you need: Few-shot rl fine-tuning of language models. *ArXiv preprint*, abs/2506.06395, 2025a.

Zongxia Li, Yapei Chang, Yuhang Zhou, Xiyang Wu, Zichao Liang, Yoo Yeon Sung, and Jordan Lee Boyd-Graber. Semantically-aware rewards for open-ended r1 training in free-form generation. *ArXiv preprint*, abs/2506.15068, 2025b.

Xiao Liang, Zhongzhi Li, Yeyun Gong, Yelong Shen, Ying Nian Wu, Zhijiang Guo, and Weizhu Chen. Beyond pass@ 1: Self-play with variational problem synthesis sustains rlvr. *arXiv preprint arXiv:2508.14029*, 2025.

Zi Lin, Sheng Shen, Jingbo Shang, Jason Weston, and Yixin Nie. Learning to solve and verify: A self-play framework for code and test generation. *ArXiv preprint*, abs/2502.14948, 2025.

Bo Liu, Chuanyang Jin, Seungone Kim, Weizhe Yuan, Wenting Zhao, Ilia Kulikov, Xian Li, Sainbayar Sukhbaatar, Jack Lanchantin, and Jason Weston. Spice: Self-play in corpus environments improves reasoning. *arXiv preprint arXiv:2510.24684*, 2025a.

Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025b.

Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, et al. General-reasoner: Advancing llm reasoning across all domains. *ArXiv preprint*, abs/2505.14652, 2025.

John X Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G Edward Suh, Alexander M Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. How much do language models memorize? *arXiv preprint arXiv:2505.24832*, 2025.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

Julien Pourcel, Cédric Colas, and Pierre-Yves Oudeyer. Self-improving language models for evolutionary program synthesis: A case study on arc-agi. 2025.

Mihir Prabhudesai, Lili Chen, Alex Ippoliti, Katerina Fragkiadaki, Hao Liu, et al. Maximizing confidence alone improves reasoning. *ArXiv preprint*, abs/2505.22660, 2025.

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. In *First Conference on Language Modeling (COLM)*.

Sheikh Shafayat, Fahim Tajwar, Ruslan Salakhutdinov, Jeff Schneider, and Andrea Zanette. Can large reasoning models self-train? *ArXiv preprint*, abs/2505.21444, 2025.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. Efficient reinforcement finetuning via adaptive curriculum learning. *ArXiv preprint*, abs/2504.05520, 2025.

Mehrangiz shoaa kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, et al. Big-bench extra hard. In *Annual Meeting of the Association for Computational Linguistics*, 2025.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017a.

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, et al. Mastering the game of go without human knowledge. *Nature*, 550, 2017b.

Yi Su, Dian Yu, Linfeng Song, Juntao Li, Haitao Mi, Zhaopeng Tu, Min Zhang, and Dong Yu. Crossing the reward bridge: Expanding rl with verifiable rewards across diverse domains. *arXiv preprint arXiv:2503.23829*, 2025.

Sainbayar Sukhbaatar, Zeming Lin, Ilya Kostrikov, Gabriel Synnaeve, Arthur Szlam, and Rob Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. *arXiv preprint arXiv:1703.05407*, 2017.

Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, et al. A survey on self-evolution of large language models. *ArXiv preprint*, abs/2404.14387, 2024.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.

Yinjie Wang, Ling Yang, Ye Tian, Ke Shen, and Mengdi Wang. Co-evolving llm coder and unit tester via reinforcement learning. *ArXiv preprint*, abs/2506.03136, 2025.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025a.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, et al. Qwen3 technical report. *ArXiv preprint*, abs/2505.09388, 2025b.

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

Kongcheng Zhang, Qi Yao, Shunyu Liu, Yingjie Wang, Baisheng Lai, et al. Consistent paths lead to truth: Self-rewarding reinforcement learning for llm reasoning. *ArXiv preprint*, abs/2506.08745, 2025.

Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin Xu, Matthieu Lin, Shenzhi Wang, Qingyun Wu, Zilong Zheng, and Gao Huang. Absolute zero: Reinforced self-play reasoning with zero data. *arXiv preprint arXiv:2505.03335*, 2025a.

Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Xiaodong Song. Learning to reason without external rewards. *ArXiv preprint*, abs/2505.19590, 2025b.

Yulai Zhao, Haolin Liu, Dian Yu, S. Y. Kung, Haitao Mi, and Dong Yu. One token to fool llm-as-a-judge. volume abs/2507.08794, 2025c.

Yaowei Zheng, Junting Lu, Shenzhi Wang, Zhangchi Feng, Dongdong Kuang, et al. Easyr1: An efficient, scalable, multi-modality rl training framework. 2025.

Yuxin Zuo, Kaiyan Zhang, Shang Qu, Li Sheng, Xuekai Zhu, et al. Ttrl: Test-time reinforcement learning. *ArXiv preprint*, abs/2504.16084, 2025.

# A Appendix

## A.1 Training Details

Our implementation is based on the EasyR1 (Zheng et al., 2025) and R-Zero (Huang et al., 2025) codebases. Figure 5 shows the solver training curve. The solver is trained for 100 steps, while the challenger is trained for 50 steps. For both models, we use a batch size of 512, a rollout number of 8, and a learning rate of 5e-7. The challenger and solver are trained iteratively, with the challenger updated for 5 steps followed by 10 steps of solver training. We do not alternate on every step because a single solver update does not significantly change its ability; performing several solver updates per iteration improves training efficiency. At the beginning of each cycle (steps 11, 21, 31, ...), the solver is trained on newly generated questions from the challenger. As the challenger generates harder questions,



Figure 5: Training curve of the solver (Qwen3-8B-Base), trained for 100 steps while alternating with the challenger.

the solver's accuracy reward drops sharply at each refresh cycle. For solver training, the reward weights are set to 0.1 for format correctness and 0.9 for accuracy. For challenger training, the uncertainty reward is computed by sampling eight responses from the current solver.
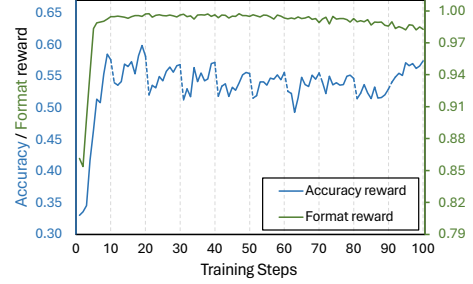
## A.2 Training Hyperparameter

This section summarizes the hyperparameters for the Solver and Challenger training stages. All experiments were conducted using BFloat16 (BF16) mixed precision and FlashAttention 2.

### A.2.1 Solver Training

- **Global Batch Size**: 512

- **Learning Rate**: $5 \times 10^{-7}$

- **KL Penalty Coefficient** ($\lambda_{KL}$): $1 \times 10^{-2}$

- **Max Steps**: 100

- **Number of Rollouts**: 8

- **Rollout Temperature**: 1.0

- **Rollout Top-p**: 0.99

### A.2.2 Challenger Training

- **Global Batch Size**: 512

- **Learning Rate**: $5 \times 10^{-7}$

- **KL Penalty Coefficient** ($\lambda_{KL}$): $1 \times 10^{-2}$

- **Max Steps**: 50

- **Number of Rollouts**: 8

- **Rollout Temperature**: 1.0

- **Rollout Top-p**: 0.99

### A.3  Prompt Templates

This section presents the exact prompt templates used for the solver and challenger models.

All highlighted items marked as "<xxx>" will be replaced with the actual input.

---

**Solver Prompt Template**

**System Message:**
Please reason step by step, and put your final answer within \boxed{}.
**User Message:**
`<problem_statement>`

---

**Challenger Prompt Template**

**System Message:**
You are an expert problem setter. First, in your private scratchpad, carefully design a brand-new, non-trivial reasoning problem. The subject may be drawn from any high-school or university discipline (mathematics, science, history, literature, social studies, etc.). The problem must be challenging enough that fewer than 30% of college students would be able to solve it correctly. You will be provided with some example questions, but you are encouraged to brainstorm freely and not be limited by those examples. Finally, output exactly in the following format:

{question}
{The complete problem statement on one or more lines}
{/question}

**User Message:**
`<Example 1>`
`<Example 2>`
`<Example 3>`
`<Example 4>`
`<Example 5>`
Generate a brand-new, challenging reasoning questions now. Each must follow the required format within {question} {/question}.

---

### A.4  GPT-4o Judge Prompt

To programmatically evaluate the correctness of answers on mathematical benchmarks where the final answer can be complex (e.g., simplified expressions), we use GPT-4o as a judge. The exact prompt and configuration used for this evaluation are detailed below.

---

### Configuration for GPT-4o as Judge

- **Model**: `gpt-4o`

- **Temperature**: 0.0

**System Message:**

You are a helpful assistant.

**User Message Template:**

Look at the following two expressions (answers to a math problem) and judge whether they are equivalent. Only perform trivial simplifications

Examples:

Expression 1: $2x + 3$

Expression 2: $3 + 2x$

Yes

Expression 1: 3/2

Expression 2: 1.5

Yes

Expression 1: $x^2 + 2x + 1$

Expression 2: $y^2 + 2y + 1$

No

Expression 1: $x^2 + 2x + 1$

Expression 2: $(x + 1)^2$

Yes

Expression 1: 3245/5

Expression 2: 649

No

(these are actually equal, don't mark them equivalent if you need to do nontrivial simplifications)

Expression 1: 2/(-3)

Expression 2: -2/3

Yes

(trivial simplifications are allowed)

Expression 1: 72 degrees

Expression 2: 72

Yes

(give benefit of the doubt to units)

Expression 1: 64

Expression 2: 64 square feet

Yes

(give benefit of the doubt to units)

—

YOUR TASK

Respond with only "Yes" or "No" (without quotes). Do not include a rationale.

Expression 1: `<Expression 1>`

Expression 2: `<Expression 2>`