

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

Abstract

We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-tuning (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1-Zero, DeepSeek-R1, and six dense models (1.5B, 7B, 8B, 14B, 32B, 70B) distilled from DeepSeek-R1 based on Qwen and Llama.

DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

Abstract

*f J i PDF1 GPT-Academic • y i z V2 i z V2< +Latex i z V2 V2 HCR Y t \ @ i z V2 V2 V2 V2 i z V2
œ ÷ i z V2 + v i z V2 : z V2 y i z V2 GitHub O @ https://github.com/binary-husky/gpt_academic/ SM i z V2 :
qwen-plus SM i z V2) i z V2 3.0.2 : t 2 b i z V2 , , i z V2 i c q z V2 1/2 b i z V2 d f J

[illegible]

Figure 1 j Benchmark performance of DeepSeek-R1.

Figure 1 j DeepSeek-R1, Ĩ Ĩ 1/2 Ĩ 1/2

Contents

1	Introduction	4
1.1	Contributions	5
1.2	Summary of Evaluation Results	6
2	Approach	6
2.1	Overview	6
2.2	DeepSeek-R1-Zero: Reinforcement Learning on the Base Model	7
2.2.1	Reinforcement Learning Algorithm	7
2.2.2	Reward Modeling	7
2.2.3	Training Template	8
2.2.4	Performance, Self-evolution Process and Aha Moment of DeepSeek-R1-Zero	8
2.3	DeepSeek-R1: Reinforcement Learning with Cold Start	11
2.3.1	Cold Start	12
2.3.2	Reasoning-oriented Reinforcement Learning	12
2.3.3	Rejection Sampling and Supervised Fine-Tuning	13
2.3.4	Reinforcement Learning for all Scenarios	13
2.4	Distillation: Empower Small Models with Reasoning Capability	14

Contents

1	Introduction	4
1.1	Contributions	5
1.2	Summary of Evaluation Results	5
2	Approach	6
2.1	Overview	6
2.2	DeepSeek-R1-Zero: Ĩ Ĩ 1/2 Ĩ 1/2 : f`	6
2.2.1	Reinforcement Learning Algorithm	6
2.2.2	Reward Modeling	7
2.2.3	Training Template	7
2.2.4	Performance, Self-evolution Process and Aha Moment of DeepSeek-R1-Zero	7
2.3	DeepSeek-R1: Reinforcement Learning with Cold Start	10
2.3.1	Cold Start	10
2.3.2	Reasoning-oriented Reinforcement Learning	11
2.3.3	Rejection Sampling and Supervised Fine-Tuning	11
2.3.4	Reinforcement Learning for all Scenarios	12
2.4	Distillation: Empower Small Models with Reasoning Capability	12

3	Experiment	14
3.1	DeepSeek-R1 Evaluation	16
3.2	Distilled Model Evaluation	17
4	Discussion	18
4.1	Distillation v.s. Reinforcement Learning	18
4.2	Unsuccessful Attempts	18
5	Conclusion, Limitations, and Future Work	20
A	Contributions and Acknowledgments	25

3	Experiment	12
3.1	DeepSeek-R1 Evaluation	13
3.2	Distilled Model Evaluation	15
4	Discussion	15
4.1	Distillation v.s. Reinforcement Learning	15
4.2	Unsuccessful Attempts	16
5	Conclusion, Limitations, and Future Work	17
A	Contributions and Acknowledgments	22

into account prompts from all scenarios. After these steps, we obtained a checkpoint referred to as DeepSeek-R1, which achieves performance on par with OpenAI-o1-1217.

We further explore distillation from DeepSeek-R1 to smaller dense models. Using Qwen2.5-72B (Qwen, 2024b) as the base model, direct distillation from DeepSeek-R1 outperforms applying RL on it. This demonstrates that the reasoning patterns discovered by larger base models are crucial for improving reasoning capabilities. We open-source the distilled Qwen and Llama (Dubey et al., 2024) series. Notably, our distilled 72B model outperforms state-of-the-art open-source QwQ-72B-Preview (Qwen, 2024a) by a large margin, and the distilled 72B and 144B models set a new record on the reasoning benchmarks among dense models.

1.1. Contributions

Post-Training: Large-Scale Reinforcement Learning on the Base Model

- We directly apply RL to the base model without relying on supervised fine-tuning (SFT) as a preliminary step. This approach allows the model to explore chain-of-thought (CoT) for solving complex problems, resulting in the development of DeepSeek-R1-Zero. DeepSeek-R1-Zero demonstrates capabilities such as self-verification, reflection, and generating long CoTs, marking a significant milestone for the research community. Notably, it is the first open research to validate that reasoning capabilities of LLMs can be incentivized purely through RL, without the need for SFT. This breakthrough paves the way for future advancements in this area.
- We introduce our pipeline to develop DeepSeek-R1. The pipeline incorporates two RL stages aimed at discovering improved reasoning patterns and aligning with human preferences, as well as two SFT stages that serve as the seed for the model's reasoning and non-reasoning capabilities. We believe the pipeline will benefit the industry by creating better models.

Distillation: Smaller Models Can Be Powerful Too

- We demonstrate that the reasoning patterns of larger models can be distilled into smaller models, resulting in better performance compared to the reasoning patterns discovered through RL on small models. The open source DeepSeek-R1, as well as its API, will benefit the research community to distill better smaller models in the future.
- Using the reasoning data generated by DeepSeek-R1, we fine-tuned several dense models that are widely used in the research community. The evaluation results demonstrate that the distilled smaller dense models perform exceptionally well on benchmarks. DeepSeek-R1-Distill-Qwen-7B achieves 55.5% on AIME 2024, surpassing QwQ-32B-Preview. Additionally, DeepSeek-R1-Distill-Qwen-32B scores 72.6% on AIME 2024, 94.3% on MATH-500,

1.1. Contributions

$$- \frac{1}{2} : (i \otimes \frac{1}{2} \prec \frac{1}{2} \prec \frac{1}{2} : \otimes f \prec$$

- [illegible]

$$s \cdot \frac{1}{2} \quad " \quad ! \quad < \quad - \quad i \quad \frac{1}{2} \quad \frac{1}{2}$$

- $\tilde{t} \in \frac{1}{2} \mathbb{N}_2 f' ! < _, _ ! \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \cdot \text{Of} _, ! < - \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \tilde{t} \mathbb{N}_2 \mathbb{N}_2 ! < \tilde{t} \in \frac{1}{2} \mathbb{N}_2 ! w \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \cdot _, \text{DeepSeek-R1} \tilde{t} \in \mathbb{N}_2 \text{API} \odot \tilde{Z} \vee >: (* e _, \cdot \tilde{t} \in \frac{1}{2} \mathbb{N}_2 < ! <$
- • (1 DeepSeek-R1 $_, _ p n \tilde{t} \in \mathbb{N}_2 \dagger \tilde{t} \in \mathbb{N}_2 \vee >: - \cdot \tilde{t} \in \mathbb{N}_2 _, \tilde{t} \in \frac{1}{2} \mathbb{N}_2 ! < \tilde{t} \in \mathbb{N}_2 \odot \mathbb{N}_2 \mathbb{N}_2: \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \cdot _, f \tilde{t} \in \frac{1}{2} \mathbb{N}_2 (\tilde{t} \in \frac{1}{2} \mathbb{N}_2 \in \mathbb{N}_2^\circ \tilde{t} \in \frac{1}{2} \mathbb{N}_2 \text{DeepSeek-R1-Distill-Qwen-7B} (\text{AIME 2024} \tilde{t} \in \mathbb{N}_2 55.5\% \dots \tilde{t} \in \mathbb{N}_2 \text{QwQ-32B-Preview} d \text{DeepSeek-R1-Distill-Qwen-32B} (\text{AIME 2024} - 72.6\% (\text{MATH-500} - 94.3\% (\text{LiveCodeBench} - 57.2\% \tilde{t} \in \mathbb{N}_2 \odot \mathbb{N}_2 W \tilde{Z} K M _, \cdot ! < \vee o1-mini \tilde{t} \in \frac{1}{2} \mathbb{N}_2 S \tilde{t} \in \frac{1}{2} \mathbb{N}_2: \cdot \dagger \tilde{t} \in \mathbb{N}_2 \text{Qwen2.5} \mathbb{N}_2 \text{Llama3} \tilde{t} \in \frac{1}{2} \mathbb{N}_2 1.5B \ 7B \ 8B \ 14B \ 32B \mathbb{N}_2 70B \tilde{t} \in \mathbb{N}_2 \mathbb{N}_2 \cdot ! < \tilde{t} \in \mathbb{N}_2$

1.2. Summary of Evaluation Results

- DeepSeek-R1 (AIME 2024) Pass@1 79.8%, OpenAI-o1-1217 (MATH-500) 97.3%
OpenAI-o1-1217 vs DeepSeek-R1 U.S. Codeforces Elo 2,029 vs 1,963
DeepSeek-R1 vs DeepSeek-V3 vs GPT-4o vs Gemini-1.5 Pro vs Claude-3.5 Sonnet vs LLaMA-3.3 70B Instruct
X (Elo) vs Ho
- MMLU-Pro GPQA Diamond DeepSeek-R1 vs DeepSeek-V3 (MMLU: 90.8% vs 84.0%, GPQA Diamond: 71.5% vs 69.1%)
OpenAI-o1-1217 vs DeepSeek-R1 vs Gemini-1.5 Pro vs Claude-3.5 Sonnet vs LLaMA-3.3 70B Instruct
SimpleQA DeepSeek-R1 vs DeepSeek-V3 vs Gemini-1.5 Pro vs Claude-3.5 Sonnet vs LLaMA-3.3 70B Instruct

and 57.2% on LiveCodeBench. These results significantly outperform previous open-source models and are comparable to o1-mini. We open-source distilled 1.5B, 7B, 8B, 14B, 32B, and 70B checkpoints based on Qwen2.5 and Llama3 series to the community.

1.2. Summary of Evaluation Results

- **Reasoning tasks:** (1) DeepSeek-R1 achieves a score of 79.8% Pass@1 on AIME 2024, slightly surpassing OpenAI-o1-1217. On MATH-500, it attains an impressive score of 97.3%, performing on par with OpenAI-o1-1217 and significantly outperforming other models. (2) On coding-related tasks, DeepSeek-R1 demonstrates expert level in code competition tasks, as it achieves 2,029 Elo rating on Codeforces outperforming 96.3% human participants in the competition. For engineering-related tasks, DeepSeek-R1 performs slightly better than DeepSeek-V3, which could help developers in real world tasks.
- **Knowledge:** On benchmarks such as MMLU, MMLU-Pro, and GPQA Diamond, DeepSeek-R1 achieves outstanding results, significantly outperforming DeepSeek-V3 with scores of 90.8% on MMLU, 84.0% on MMLU-Pro, and 71.5% on GPQA Diamond. While its performance is slightly below that of OpenAI-o1-1217 on these benchmarks, DeepSeek-R1 surpasses other closed-source models, demonstrating its competitive edge in educational tasks. On the factual benchmark SimpleQA, DeepSeek-R1 outperforms DeepSeek-V3, demonstrating its capability in handling fact-based queries. A similar trend is observed where OpenAI-o1 surpasses 4o on this benchmark.
- **Others:** DeepSeek-R1 also excels in a wide range of tasks, including creative writing, general question answering, editing, summarization, and more. It achieves an impressive length-controlled win-rate of 87.6% on AlpacaEval 2.0 and a win-rate of 92.3% on ArenaHard, showcasing its strong ability to intelligently handle non-exam-oriented queries. Additionally, DeepSeek-R1 demonstrates outstanding performance on tasks requiring long-context understanding, substantially outperforming DeepSeek-V3 on long-context benchmarks.

2. Approach

2.1. Overview

Previous work has heavily relied on large amounts of supervised data to enhance model performance. In this study, we demonstrate that reasoning capabilities can be significantly improved through large-scale reinforcement learning (RL), even without using supervised fine-tuning (SFT) as a cold start. Furthermore, performance can be further enhanced with the inclusion of a small amount of cold-start data. In the following sections, we present: (1) DeepSeek-R1-Zero, which applies RL directly to the base model without any SFT data, and

• v i 占 DeepSeek-R1 (• i 占 ½ h ° i 占 ½ j ½ T M \ (i 占 ½ j ½ ' ;
i 占 ½ f (AlpacaEval 2.0 ž ° t i 占 ½ a i 占 ½ • | S 6 i 占 ½ 77.6% (ArenaHard
„ i 占 ½ 92.3% U: t v z i 占 ½ ^ i 占 ½ j ½ j ½ i 占 ½ d DeepSeek-R1
(• • ‡ j ½ i 占 ½ h ° i 占 ½ > W Ž DeepSeek-V3 (• ‡ i 占 ½ j ½
„ h °

2. Approach

2.1. Overview

HM, vi 8% ½V' ½gþne 8! < ' i ½, v- i ½ ½ i ½ ½
! : f` RL s• • (i 8 SFT \ : • / " " i ½ i ½ ½ ½ d
e ½ / " pñi ½ ½ ½ 8! i ½ i ½ i ½ i ½ i ½ ½ (1) DeepSeek-R1-
Zero f i ½ RL" (Ž i ½ < i ½ ½ SFT pñ (2) DeepSeek-R1 f i ½ p C*
• i ½ ½ ½ : < ® „ ½ ½ ½ ½ (RL 3) i DeepSeek-R1, • " i ½ ½ < i ½ ½
! <

2.2. DeepSeek-R1-Zero: $\tilde{\mathcal{L}}_{V_2@V_2} \ll \mathcal{L}_{V_2@V_2}$

: f` (" i j ½ h° i z W, Hoe c, i j M K i x @ i j ½ E 7 (Shao et al., 2024; Wang et al., 2023) 6 i j ½ x % j ½ V Z i g p n 6 i j ½ p n ^ 8 i j ½ , - i j ½ ' < i j ½ < LLMs (i j ½ g p n , j ½ j ½ j ½ i j ½ \> j ½ s i f j ½ j ½ j ½ j ½ f ` i j ½ ° , i j ½ i j ½ i j ½ H ∈ • i j ½ j ½ j ½ R — i j ½ U : > ä ° t K , i j ½ j ½ i j ½ j ½ : j ½ j ½ ÷ < , j ½ j ½ j ½

2.2.1. Reinforcement Learning Algorithm

$\tilde{V}_2 \tilde{V}_2$: \dagger , \vdash $f \sim -i \tilde{V}_2$ $\tilde{V}_2 \dagger \tilde{V}_2$ GRPO (Shao et al., 2024) $\tilde{V}_2 \tilde{V}_2 \dagger 8$ $Ve! < \tilde{V}_2 \tilde{A}^\circ \nabla! < \nabla \tilde{V}_2 - O_i \tilde{V}_2 w$
 $Se \tilde{V}_2 \tilde{V}_2 \tilde{V}_2 \tilde{V}_2 \tilde{V}_2$ GRPO $\tilde{V}_2 c_{\lambda=3} - \tilde{V}_2 \tilde{V}_2 \tilde{V}_2 \tilde{V}_2 = \tilde{V}_2 = 2$, $\vdash = g \ 6$ \tilde{V}_2
 $\tilde{V}_2 \tilde{V}_2 \tilde{V}_2$ $Ve! < c_\lambda$

$$\min_{j=1}^J \frac{1}{J} \sum_{j=1}^J \min_{\mathbf{c}_{j,1}, \mathbf{c}_{j,2}} \frac{\|\mathbf{c}_{j,1} - \mathbf{c}_{j,2}\|_2}{\|\mathbf{c}_{j,1}\|_2} \quad (1)$$

$$D_{\parallel} = \frac{C_{\parallel} C_{45}}{C_{\parallel} C_{45}} = \frac{C_{45} \sqrt{z_j} \sqrt{z_j}}{C_{\parallel} \sqrt{z_j} \sqrt{z_j}} \log \frac{C_{45} \sqrt{z_j} \sqrt{z_j}}{C_{\parallel} \sqrt{z_j} \sqrt{z_j}} \quad (2)$$

(2) DeepSeek-R1, which applies RL starting from a checkpoint fine-tuned with thousands of long Chain-of-Thought (CoT) examples. 3) Distill the reasoning capability from DeepSeek-R1 to small dense models.

2.2. DeepSeek-R1-Zero: Reinforcement Learning on the Base Model

Reinforcement learning has demonstrated significant effectiveness in reasoning tasks, as evidenced by our previous works (Shao et al., 2024; Wang et al., 2023). However, these works heavily depended on supervised data, which are time-intensive to gather. In this section, we explore the potential of LLMs to develop reasoning capabilities **without any supervised data**, focusing on their self-evolution through a pure reinforcement learning process. We start with a brief overview of our RL algorithm, followed by the presentation of some exciting results, and hope this provides the community with valuable insights.

2.2.1. Reinforcement Learning Algorithm

Group Relative Policy Optimization In order to save the training costs of RL, we adopt Group Relative Policy Optimization (GRPO) (Shao et al., 2024), which foregoes the critic model that is typically the same size as the policy model, and estimates the baseline from group scores instead. Specifically, for each question q , GRPO samples a group of outputs $f_{1,2,\dots,g}$ from the old policy $\pi_{\theta_{old}}$ and then optimizes the policy model π_{θ} by maximizing the following objective:

$$J(\theta) = \mathbb{E}_{q \sim \mathcal{D}} \sum_{j=1}^g \pi_{\theta}(f_j|q) \left[\frac{1}{g} \sum_{i=1}^g \min \left(\frac{\pi_{\theta}(f_j|q)}{\pi_{\theta_{old}}(f_j|q)} \right)^{\gamma}, \text{clip} \left(\frac{\pi_{\theta}(f_j|q)}{\pi_{\theta_{old}}(f_j|q)}, 1-\gamma, 1+\gamma \right) \right] V(f_j), \quad (1)$$

$$D_{KL}(\pi_{\theta} \parallel \pi_{\theta_{old}}) = \frac{C_{\theta_{old}}(f_j|q)}{C_{\theta}(f_j|q)} \log \frac{C_{\theta_{old}}(f_j|q)}{C_{\theta}(f_j|q)}, \quad (2)$$

where γ and V are hyper-parameters, and γ is the advantage, computed using a group of rewards $f_{\theta_1, \theta_2, \dots, \theta_g}$ corresponding to the outputs within each group:

$$\gamma = \frac{\theta_1 - m40\langle f_{\theta_1, \theta_2}, \dots, \theta_g \rangle}{SB3\langle f_{\theta_1, \theta_2}, \dots, \theta_g \rangle}. \quad (3)$$

2.2.2. Reward Modeling

The reward is the source of the training signal, which decides the optimization direction of RL. To train DeepSeek-R1-Zero, we adopt a rule-based reward system that mainly consists of two types of rewards:

- **Accuracy rewards:** The accuracy reward model evaluates whether the response is correct.

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within `<think>` `</think>` and `<answer>` `</answer>` tags, respectively, i.e., `<think>` reasoning process here `</think>` `<answer>` answer here `</answer>`. User: `prompt`. Assistant:

Table 1 j DeepSeek-R1-Zero “ ! • textcolorredprompt (- i 2 2 2 wS “ “ i 2 2 2

$$v = \gamma \mathbb{E}_{V \sim \dots} \left[\frac{1}{g} \sum_{i=1}^g \left(\frac{\pi_{\theta}(f_i|q)}{\pi_{\theta_{old}}(f_i|q)} \right)^{\gamma} \right] f_{\theta_1, \theta_2, \dots, \theta_g} \text{ — } i \text{ 2 } \frac{1}{2}$$
$$\gamma = \frac{\theta_1 - m40\langle f_{\theta_1, \theta_2}, \dots, \theta_g \rangle}{SB3\langle f_{\theta_1, \theta_2}, \dots, \theta_g \rangle}. \quad (3)$$

2.2.2. Reward Modeling

$V \pm / - i \text{ 2 } “ e \bullet i \text{ 2 } RL “ i \text{ 2 } 2: \dagger - i \text{ 2 } \text{DeepSeek-R1-Zero} i \text{ 2 } 2 \dagger \ast$
 $i \text{ 2 } 2 \text{ 2 } “ V \pm i \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } 1 \text{ 2 } i \text{ 2 } 2 “ V \pm i \text{ 2 } 2$

- $i \text{ 2 } “ V \pm i \text{ 2 } 2 V \pm ! < i \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } \& c n < , (p f i \text{ 2 } w n s ' i \text{ 2 } 2 \text{ 2 } 2$
 $i \text{ 2 } 2 ! < \bullet i \text{ 2 } 2 < i \text{ 2 } 2 i \text{ 2 } 2 < , (i \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2$
 $“ c n ' \text{ 2 } i \text{ 2 } 2 7 i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 i \text{ 2 } 2 i \text{ 2 } 2 n “ s i “ K i \text{ 2 } 2 i \text{ 2 } 2$
- $< V \pm d \dagger i \text{ 2 } 2 V \pm ! < i \text{ 2 } 2 \text{ 2 } 2 \dagger i \text{ 2 } 2 V \pm ! < i \text{ 2 } 2: 6 ! <$
 $“ i \text{ 2 } 2 (' < think > ' \text{ 2 } < / think > ' \sim K i \text{ 2 } 2$
- $i \text{ 2 } 2 i \text{ 2 } \text{DeepSeek-R1-Zero} i \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 < i \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 <$
 $(' i \text{ 2 } 2: f \text{ 2 } i \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } \times V \pm \ll i \text{ 2 } 2 V i \text{ 2 } 2 - i \text{ 2 } 2 < \bullet \bullet “ - i \text{ 2 } 2$
- $v \bullet t \ast - i \text{ 2 } 2 B$

2.2.3. Training Template

$: \dagger - i \text{ 2 } \text{DeepSeek-R1-Zero} i \text{ 2 } 2 i \text{ 2 } 2 \ast \in U “ ! \bullet i \text{ 2 } 2 \text{ 2 } 2 < u \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } “$
 $i \text{ 2 } 2 h 1 @: i \text{ 2 } 2 \bullet B \text{DeepSeek-R1-Zero} - H \ast “ i \text{ 2 } 2 6 / i \text{ 2 } 2$
 $i \text{ 2 } 2 i \text{ 2 } _ P 6 (i \text{ 2 } 2 \text{ 2 } 2 \bullet M i \text{ 2 } 2 i \text{ 2 } 2 “ O i \text{ 2 } 2 < , : 6 i \text{ 2 } 2 “$
 $“ \bullet y s “ i \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } 2 (R L i \text{ 2 } 2 , i \text{ 2 } 2 i \text{ 2 } 2$

2.2.4. Performance, Self-evolution Process and Aha Moment of DeepSeek-R1-Zero

Performance of DeepSeek-R1-Zero $i \text{ 2 } 2$: $\dagger \text{DeepSeek-R1-Zero} (\text{AIME 2024} i \text{ 2 } 2 \text{ 2 } 2$
 $(t \ast RL - i \text{ 2 } 2 \text{ 2 } 2 “ i \text{ 2 } 2 \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } 2 \text{ 2 } \bullet @ RL - i \text{ 2 } 2 i \text{ 2 } 2 \text{ 2 } \text{DeepSeek-R1-Zero} h \circ$
 $i \text{ 2 } 2 \text{ 2 } i \text{ 2 } 2 ' i \text{ 2 } 2 < \text{—} i \text{ 2 } 2 / \text{AIME 2024} “ s G \text{ pass@1 } p > W i \text{ 2 } 2 i \text{ 2 } 2 “$
 $15.6 \% i \text{ 2 } 2 O i \text{ 2 } 2 \text{ 2 } a i \text{ 2 } 2 71.0 \% i \text{ 2 } 2 \text{ 2 } \text{OpenAI-o1-0912} i \text{ 2 } 2 i \text{ 2 } 2 i \text{ 2 } 2 W 9$

Table 1 j Template for DeepSeek-R1-Zero. **prompt** will be replaced with the specific reasoning question during training.

Table 2 j Comparison of DeepSeek-R1-Zero and OpenAI o1 models on reasoning-related benchmarks.

- **Format rewards:** In addition to the accuracy reward model, we employ a format reward model that enforces the model to put its thinking process between '`<think>`' and '`</think>`' tags.

2.2.3. Training Template

such as mandating reflective reasoning or promoting particular problem-solving strategies to ensure that we can accurately observe the model’s natural progression during the RL process.

2.2.4. Performance, Self-evolution Process and Aha Moment of DeepSeek-R1-Zero

Table 2 j DeepSeek-R1-Zero & OpenAI o1! < (" ï § V₂ L K₂V₂L ½L f₂

$$\ddot{r} = \frac{1}{2} + \ddot{r}_{\text{RL}} - \ddot{\epsilon} \frac{1}{2} ! < ' \ddot{r} \approx \frac{1}{2} H'$$
[illegible]

Figure 3 j The average response length of DeepSeek-R1-Zero on the training set during the RL process. DeepSeek-R1-Zero naturally learns to solve reasoning tasks with more thinking time.

progression without the influence of the supervised fine-tuning stage. This approach provides a clear view of how the model evolves over time, particularly in terms of its ability to handle complex reasoning tasks.

As depicted in Figure 3, the thinking time of DeepSeek-R1-Zero shows consistent improvement throughout the training process. This improvement is not the result of external adjustments but rather an intrinsic development within the model. DeepSeek-R1-Zero naturally acquires the ability to solve increasingly complex reasoning tasks by leveraging extended test-time computation. This computation ranges from generating hundreds to thousands of reasoning tokens, allowing the model to explore and refine its thought processes in greater depth.

One of the most remarkable aspects of this self-evolution is the emergence of sophisticated behaviors as the test-time computation increases. Behaviors such as reflection where the model revisits and reevaluates its previous steps and the exploration of alternative approaches to problem-solving arise spontaneously. These behaviors are not explicitly programmed but instead emerge as a result of the model's interaction with the reinforcement learning environment. This spontaneous development significantly enhances DeepSeek-R1-Zero's reasoning capabilities, enabling it to tackle more challenging tasks with greater efficiency and accuracy.

Aha Moment of DeepSeek-R1-Zero A particularly intriguing phenomenon observed during the training of DeepSeek-R1-Zero is the occurrence of an "aha moment". This moment, as illustrated in Table 3, occurs in an intermediate version of the model. During this phase, DeepSeek-R1-Zero learns to allocate more thinking time to a problem by reevaluating its initial approach. This behavior is not only a testament to the model's growing reasoning abilities

Question: If $0 < F < 1$, then the sum of the real solutions of $\sqrt{0 - \sqrt{0 - F}} = F$ is equal to

Response: <think>
 To solve the equation $\sqrt{0 - \sqrt{0 - F}} = F$, let's start by squaring both sides:
 $(\sqrt{0 - \sqrt{0 - F}})^2 = F^2 \Rightarrow 0 - \sqrt{0 - F} = F^2$
 Rearrange to isolate the inner square root term:
 $-\sqrt{0 - F} = F^2 \Rightarrow \sqrt{0 - F} = -F^2$
 $\sqrt{0 - F} = -F^2 \Rightarrow 0 - F = (-F^2)^2 \Rightarrow -F = F^4$
 $F^4 + F = 0 \Rightarrow F(F^3 + 1) = 0$
 $F = 0$ or $F^3 + 1 = 0 \Rightarrow F^3 = -1 \Rightarrow F = -1$
 ...
 Wait, wait. Wait. That's an aha moment I can flag here.
 Let's reevaluate this step-by-step to identify if the correct sum can be
 We started with the equation:
 $\sqrt{0 - \sqrt{0 - F}} = F$
 First, let's square both sides:
 $0 - \sqrt{0 - F} = F^2 \Rightarrow \sqrt{0 - F} = -F^2$
 Next, I could square both sides again, treating the equation:
 ...

Table 3 j * E., DeepSeek-R1-Zero - i H/2, J i E 1/2 1/2 i E 1/2 F + (E 1/2 E 1/2 1/2 i E 1/2 1/2 i E 1/2 1/2 i E 1/2 1/2 * J i E 1/2 1/2 • i E 1/2 1/2 1/2 F ^ , > E 1/2 Z =

i E 1/2 i E 1/2 1/2 / " p n , : f ^ i E 1/2 1/2

2.3. DeepSeek-R1: Reinforcement Learning with Cold Start

× ODeepSeek-R1-Zero i E 1/2 i E 1/2 1/2 / i E 1/2 * i E 1/2 1/2 1) i E 1/2 i E 1/2 1/2 1/2 p n \ : . / " i E 1/2 1/2 1/2 1/2 i E 1/2 1/2 6 [2) , U - i E 1/2 (7 i E 1/2 ! < • v i E 1/2 1/2 p i E 1/2 1/2 i E 1/2 1/2 CoT i E 1/2 1/2 i E 1/2 1/2 (i E 1/2 : + i E 1/2 1/2 1/2 1/2 i E 1/2 1/2 * - i DeepSeek-R1, A 4 i E 1/2 1/2 + i E 1/2 1/2 w S ,

2.3.1. Cold Start

Z DeepSeek-R1-Zero : + 2 b i E 1/2 1/2 i E 1/2 1/2 RL- i E 1/2 1/2 1/2 3 S , . / " 6 i E 1/2 i E 1/2 1/2 1/2 i E 1/2 1/2 1/2 1/2 CoT p n e ® ! < \ : i E 1/2 , RLg L : + 6 i E 1/2 1/2 p n i E 1/2 1/2 i E 1/2 1/2 1/2 (• i E 1/2 \ : : < , i E 1/2 , i E 1/2 i E 1/2 1/2 ! < & i E 1/2 1/2 i E 1/2 1/2 i E 1/2 1/2 6 i DeepSeek-R1-Zero, " i E 1/2 i E 1/2 1/2 i E 1/2 e i E 1/2

(, i E 1/2 - i E 1/2 1/2 Ca . / " p n i E 1/2 1/2 DeepSeek-V3-Base \ : RL, w i E 1/2 DeepSeek-R1-Zero, the advantages of cold start data include:

- i E 1/2 1/2 DeepSeek-R1-Zero, * s . P 6 / v ... i E 1/2 1/2 i E 1/2 1/2 1/2 1/2 : Omarkdown< e • i E 1/2 1/2 T H i E 1/2 7 • (i E 1/2 1/2 (: DeepSeek-R1 i E 1/2 1/2 p n i E 1/2 1/2 i E 1/2 1/2 * i E 1/2 1/2 i E 1/2 1/2 (i E 1/2 1/2 1/2 + > + * X
- v i E 1/2 1/2 E > i E 1/2 1/2 i E 1/2 1/2 1/2 (i E 1/2 1/2 i E 1/2 1/2 1/2 1/2 : |special_token| < "

Question: If $0 < F < 1$, then the sum of the real solutions of $\sqrt{0 - \sqrt{0 - F}} = F$ is equal to

Response: <think>

To solve the equation $\sqrt{0 - \sqrt{0 - F}} = F$, let's start by squaring both

$$\sqrt{0 - \sqrt{0 - F}}^2 = F^2 \Rightarrow 0 - \sqrt{0 - F} = F^2.$$

Rearrange to isolate the inner square root term:

$$\sqrt{0 - F^2} = 0 - \sqrt{0 - F} \Rightarrow \sqrt{0 - F^2} = -\sqrt{0 - F} \Rightarrow \sqrt{0 - F^2} = \sqrt{0 - F} \Rightarrow 0 - F^2 = 0 - F \Rightarrow F^2 = F \Rightarrow F^2 - F = 0 \Rightarrow F(F - 1) = 0 \Rightarrow F = 0 \text{ or } F = 1.$$

...

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

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

We started with the equation:

$$\sqrt{0 - \sqrt{0 - F}} = F$$

First, let's square both sides:

$$0 - \sqrt{0 - F} = F^2 \Rightarrow \sqrt{0 - F} = 0 - F^2$$

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

...

Table 3 j An interesting “aha moment” of an intermediate version of DeepSeek-R1-Zero. The model learns to rethink using an anthropomorphic tone. This is also an aha moment for us, allowing us to witness the power and beauty of reinforcement learning.

but also a captivating example of how reinforcement learning can lead to unexpected and sophisticated outcomes.

This moment is not only an “aha moment” for the model but also for the researchers observing its behavior. It underscores the power and beauty of reinforcement learning: rather than explicitly teaching the model on how to solve a problem, we simply provide it with the right incentives, and it autonomously develops advanced problem-solving strategies. The “aha moment” serves as a powerful reminder of the potential of RL to unlock new levels of intelligence in artificial systems, paving the way for more autonomous and adaptive models in the future.

Drawback of DeepSeek-R1-Zero Although DeepSeek-R1-Zero exhibits strong reasoning capabilities and autonomously develops unexpected and powerful reasoning behaviors, it faces several issues. For instance, DeepSeek-R1-Zero struggles with challenges like poor readability, and language mixing. To make reasoning processes more readable and share them with the open community, we explore DeepSeek-R1, a method that utilizes RL with human-friendly cold-start data.

2.3. DeepSeek-R1: Reinforcement Learning with Cold Start

Inspired by the promising results of DeepSeek-R1-Zero, two natural questions arise: 1) Can reasoning performance be further improved or convergence accelerated by incorporating a small

“special_token” in the prompt? 2) Can we further improve the performance of DeepSeek-R1-Zero by incorporating a small “special_token” in the prompt?

- In the prompt, we can add a “special_token” to indicate the start of the reasoning process. For example, we can use “<think>” as the special token. This can help the model to better understand the reasoning process and improve its performance.

2.3.2. Reasoning-oriented Reinforcement Learning

In this section, we will introduce the Reasoning-oriented Reinforcement Learning (ROL) framework. The ROL framework is designed to improve the reasoning performance of the model by incorporating a small “special_token” in the prompt. The ROL framework consists of three main components: (1) Prompt Engineering, (2) Reinforcement Learning, and (3) Evaluation. In the Prompt Engineering component, we add a “special_token” to the prompt to indicate the start of the reasoning process. In the Reinforcement Learning component, we use a reward function to encourage the model to generate correct reasoning steps. In the Evaluation component, we evaluate the model's reasoning performance using a set of test cases. The ROL framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning.

2.3.3. Rejection Sampling and Supervised Fine-Tuning

In this section, we will introduce the Rejection Sampling and Supervised Fine-Tuning (RSFT) framework. The RSFT framework is designed to improve the reasoning performance of the model by incorporating a small “special_token” in the prompt. The RSFT framework consists of two main components: (1) Rejection Sampling and (2) Supervised Fine-Tuning. In the Rejection Sampling component, we use a rejection sampling algorithm to select high-quality reasoning steps from the model's output. In the Supervised Fine-Tuning component, we use the selected reasoning steps to fine-tune the model. The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning.

The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning. The RSFT framework consists of two main components: (1) Rejection Sampling and (2) Supervised Fine-Tuning. In the Rejection Sampling component, we use a rejection sampling algorithm to select high-quality reasoning steps from the model's output. In the Supervised Fine-Tuning component, we use the selected reasoning steps to fine-tune the model. The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning.

The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning. The RSFT framework consists of two main components: (1) Rejection Sampling and (2) Supervised Fine-Tuning. In the Rejection Sampling component, we use a rejection sampling algorithm to select high-quality reasoning steps from the model's output. In the Supervised Fine-Tuning component, we use the selected reasoning steps to fine-tune the model. The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning.

The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning. The RSFT framework consists of two main components: (1) Rejection Sampling and (2) Supervised Fine-Tuning. In the Rejection Sampling component, we use a rejection sampling algorithm to select high-quality reasoning steps from the model's output. In the Supervised Fine-Tuning component, we use the selected reasoning steps to fine-tune the model. The RSFT framework has been shown to significantly improve the reasoning performance of the model, especially in tasks that require complex reasoning.

amount of high-quality data as a cold start? 2) How can we train a user-friendly model that not only produces clear and coherent Chains of Thought (CoT) but also demonstrates strong general capabilities? To address these questions, we design a pipeline to train DeepSeek-R1. The pipeline consists of four stages, outlined as follows.

2.3.1. Cold Start

Unlike DeepSeek-R1-Zero, to prevent the early unstable cold start phase of RL training from the base model, for DeepSeek-R1 we construct and collect a small amount of long CoT data to fine-tune the model as the initial RL actor. To collect such data, we have explored several approaches: using few-shot prompting with a long CoT as an example, directly prompting models to generate detailed answers with reflection and verification, gathering DeepSeek-R1-Zero outputs in a readable format, and refining the results through post-processing by human annotators.

In this work, we collect thousands of cold-start data to fine-tune the DeepSeek-V3-Base as the starting point for RL. Compared to DeepSeek-R1-Zero, the advantages of cold start data include:

- **Readability:** A key limitation of DeepSeek-R1-Zero is that its content is often not suitable for reading. Responses may mix multiple languages or lack markdown formatting to highlight answers for users. In contrast, when creating cold-start data for DeepSeek-R1, we design a readable pattern that includes a summary at the end of each response and filters out responses that are not reader-friendly. Here, we define the output format as `|special_token|<reasoning_process>|special_token|<summary>`, where the reasoning process is the CoT for the query, and the summary is used to summarize the reasoning results.
- **Potential:** By carefully designing the pattern for cold-start data with human priors, we observe better performance against DeepSeek-R1-Zero. We believe the iterative training is a better way for reasoning models.

2.3.2. Reasoning-oriented Reinforcement Learning

After fine-tuning DeepSeek-V3-Base on the cold start data, we apply the same large-scale reinforcement learning training process as employed in DeepSeek-R1-Zero. This phase focuses on enhancing the model’s reasoning capabilities, particularly in reasoning-intensive tasks such as coding, mathematics, science, and logic reasoning, which involve well-defined problems with clear solutions. During the training process, we observe that CoT often exhibits language mixing, particularly when RL prompts involve multiple languages. To mitigate the issue of language mixing, we introduce a language consistency reward during RL training, which is calculated

2.3.4. Reinforcement Learning for all Scenarios

During the RL training process, we observe that CoT often exhibits language mixing, particularly when RL prompts involve multiple languages. To mitigate the issue of language mixing, we introduce a language consistency reward during RL training, which is calculated as follows: $R_{LC} = \frac{1}{L} \sum_{i=1}^L \mathbb{1}(l_i = l^*)$, where l_i is the language of the i -th token in the CoT, l^* is the language of the prompt, and $\mathbb{1}(\cdot)$ is the indicator function. The total reward is the sum of the performance reward and the language consistency reward. We use the same training process as DeepSeek-R1-Zero, but with the addition of the language consistency reward. The training process is summarized in Figure 2.3.4.

2.4. Distillation: Empower Small Models with Reasoning Capability

DeepSeek-R1-Zero has demonstrated strong performance in reasoning tasks. To empower small models with reasoning capability, we distill the knowledge from DeepSeek-R1-Zero into smaller models. We use the same distillation process as in DeepSeek-R1-Zero, but with a smaller batch size and a shorter training time. The distillation process is summarized in Figure 2.4.1.

We distill the knowledge from DeepSeek-R1-Zero into smaller models. We use the same distillation process as in DeepSeek-R1-Zero, but with a smaller batch size and a shorter training time. The distillation process is summarized in Figure 2.4.1.

We distill the knowledge from DeepSeek-R1-Zero into smaller models. We use the same distillation process as in DeepSeek-R1-Zero, but with a smaller batch size and a shorter training time. The distillation process is summarized in Figure 2.4.1.

3. Experiment

We evaluate the performance of DeepSeek-R1-Zero and DeepSeek-R1 on various reasoning tasks. The tasks include MMLU (Hendrycks et al., 2020), MMLU-Redux (Gema et al., 2024), MMLU-Pro (Wang et al., 2024), C-Eval (Huang et al., 2023), CMMLU (Li et al., 2023), IFEval (Zhou et al., 2023), FRAMES (Krishna et al., 2024), GPQA Diamond (Rein et al., 2023), SimpleQA (OpenAI, 2024c), C-SimpleQA (He et al., 2024), SWE-Bench Verified (OpenAI, 2024d), Aider¹, LiveCodeBench (Jain et al., 2024), 2024-08 – 2025-01 Codeforces², e—9K³, CNMO 2024³, AIME 2024 (MAA, 2024), d t, h°, wSei, Y, X, O, AIME 2024, MATH-500, GPQA.

¹<https://aider.chat>

²<https://codeforces.com>

³<https://www.cms.org.cn/Home/comp/comp/cid/12.html>

as the proportion of target language words in the CoT. Although ablation experiments show that such alignment results in a slight degradation in the model’s performance, this reward aligns with human preferences, making it more readable. Finally, we combine the accuracy of reasoning tasks and the reward for language consistency by directly summing them to form the final reward. We then apply RL training on the fine-tuned model until it achieves convergence on reasoning tasks.

2.3.3. Rejection Sampling and Supervised Fine-Tuning

When reasoning-oriented RL converges, we utilize the resulting checkpoint to collect SFT (Supervised Fine-Tuning) data for the subsequent round. Unlike the initial cold-start data, which primarily focuses on reasoning, this stage incorporates data from other domains to enhance the model’s capabilities in writing, role-playing, and other general-purpose tasks. Specifically, we generate the data and fine-tune the model as described below.

Reasoning data We curate reasoning prompts and generate reasoning trajectories by performing rejection sampling from the checkpoint from the above RL training. In the previous stage, we only included data that could be evaluated using rule-based rewards. However, in this stage, we expand the dataset by incorporating additional data, some of which use a generative reward model by feeding the ground-truth and model predictions into DeepSeek-V3 for judgment. Additionally, because the model output is sometimes chaotic and difficult to read, we have filtered out chain-of-thought with mixed languages, long paragraphs, and code blocks. For each prompt, we sample multiple responses and retain only the correct ones. In total, we collect about 600k reasoning related training samples.

Non-Reasoning data For non-reasoning data, such as writing, factual QA, self-cognition, and translation, we adopt the DeepSeek-V3 pipeline and reuse portions of the SFT dataset of DeepSeek-V3. For certain non-reasoning tasks, we call DeepSeek-V3 to generate a potential chain-of-thought before answering the question by prompting. However, for simpler queries, such as “hello” we do not provide a CoT in response. In the end, we collected a total of approximately 200k training samples that are unrelated to reasoning.

We fine-tune DeepSeek-V3-Base for two epochs using the above curated dataset of about 800k samples.

2.3.4. Reinforcement Learning for all Scenarios

To further align the model with human preferences, we implement a secondary reinforcement learning stage aimed at improving the model’s helpfulness and harmlessness while simultane-

Diamond Codeforces LiveCodeBench $\approx J$ $\frac{1}{2}$ $\frac{1}{2}$

[illegible][illegible][illegible]

$$\text{pass}@1 = \frac{1}{9} \sum_{l=1}^9 \frac{1}{2}$$

v- >_h: , 7 * i T 1/2cn' d i 1/2 1/2 t i 1/2 1/2' i 1/2 i 1/2 AIME 2024 i 1/2 1/2
 ¥J t q i 1/2 p • h i 1/2 Wang et al., 2022) • (64 * 7, ° : cons@64

3.1. DeepSeek-R1 Evaluation

Ĩ Ž ½Y ½Z ½Z ĩ ½İ ĩ ½Z ½Z ½Z ĩ ½ MMLU MMLU-Pro ĒGPOA Diamond DeepSeek-
 R1İ ĩ DeepSeek-V3h° ĩ ĩ ½ ½„ ' ĩ ĩ ½Z ĩ ½Z ĩ ½ RY Ž STEMİ ĩ ½Z ĩ ½„ ĩ ĩ ½Z ĩ ½
 ĩ ½ ĩ ĩ ½ ½Z: f`ž°†>W„ ĩ ĩ ½d DeepSeek-R1(FRAMES h° ĩ ĩ ½İ ĩ ½
 *•V• ‡„ ĩ ĩ ½Z ĩ ½U: †v: '„ ‡c •İ ĩ ½İ ĩ ½†“ !<(Alq“ „
 " Ēpn •İ ĩ ½„ \> (<žİ ĩ ½Z SimpleQA DeepSeek-R1...Š†DeepSeek-
 V3 ĩ ĩ ½v ĩ Ž ½Z „ ĩ ½Z ĩ ½Z ĩ ½{ <„ < ĩ >: OpenAI-o1(dİ ĩ ½Z ĩ ½İ ĩ ½
 †GPT-4o 6 (- ‡ SimpleQAİ ĩ ½Z ½Z ½DeepSeek-R1„ h° „ DeepSeek-V3 ;•

ously refining its reasoning capabilities. Specifically, we train the model using a combination of reward signals and diverse prompt distributions. For reasoning data, we adhere to the methodology outlined in DeepSeek-R1-Zero, which utilizes rule-based rewards to guide the learning process in math, code, and logical reasoning domains. For general data, we resort to reward models to capture human preferences in complex and nuanced scenarios. We build upon the DeepSeek-V3 pipeline and adopt a similar distribution of preference pairs and training prompts. For helpfulness, we focus exclusively on the final summary, ensuring that the assessment emphasizes the utility and relevance of the response to the user while minimizing interference with the underlying reasoning process. For harmlessness, we evaluate the entire response of the model, including both the reasoning process and the summary, to identify and mitigate any potential risks, biases, or harmful content that may arise during the generation process. Ultimately, the integration of reward signals and diverse data distributions enables us to train a model that excels in reasoning while prioritizing helpfulness and harmlessness.

2.4. Distillation: Empower Small Models with Reasoning Capability

To equip more efficient smaller models with reasoning capabilities like DeepSeek-R1, we directly fine-tuned open-source models like Qwen (Qwen, 2024b) and Llama (AI@Meta, 2024) using the 800k samples curated with DeepSeek-R1, as detailed in §2.3.3. Our findings indicate that this straightforward distillation method significantly enhances the reasoning abilities of smaller models. The base models we use here are Qwen2.5-Math-1.5B, Qwen2.5-Math-7B, Qwen2.5-14B, Qwen2.5-32B, Llama-3.1-8B, and Llama-3.3-70B-Instruct. We select Llama-3.3 because its reasoning capability is slightly better than that of Llama-3.1.

For distilled models, we apply only SFT and do not include an RL stage, even though incorporating RL could substantially boost model performance. Our primary goal here is to demonstrate the effectiveness of the distillation technique, leaving the exploration of the RL stage to the broader research community.

3. Experiment

Benchmarks We evaluate models on MMLU (Hendrycks et al., 2020), MMLU-Redux (Gema et al., 2024), MMLU-Pro (Wang et al., 2024), C-Eval (Huang et al., 2023), and CMMLU (Li et al., 2023), IFEval (Zhou et al., 2023), FRAMES (Krishna et al., 2024), GPQA Diamond (Rein et al., 2023), SimpleQA (OpenAI, 2024c), C-SimpleQA (He et al., 2024), SWE-Bench Verified (OpenAI, 2024d), Aider ¹, LiveCodeBench (Jain et al., 2024) (2024-08 – 2025-01), Codeforces ², Chinese National High School Mathematics Olympiad (CNMO 2024)³, and American Invitational Mathematics

¹<https://ai.der.chat>

²<https://codeforces.com>

³<https://www.cms.org.cn/Home/comp/comp/cid/12.html>

Benchmark (Metric)		Claude-3.5-Sonnet-1022	GPT-4o-0513	DeepSeek-V3	OpenAI o1-mini	OpenAI o1-1217	DeepSeek R1
	Architecture	-	-	MoE	-	-	MoE
	# Activated Params	-	-	37B	-	-	37B
	# Total Params	-	-	671B	-	-	671B
English	MMLU (Pass@1)	88.3	87.2	88.5	85.2	91.8	90.8
	MMLU-Redux (EM)	88.9	88.0	89.1	86.7	-	92.9
	MMLU-Pro (EM)	78.0	72.6	75.9	80.3	-	84.0
	DROP (3-shot F1)	88.3	83.7	91.6	83.9	90.2	92.2
	IF-Eval (Prompt Strict)	86.5	84.3	86.1	84.8	-	83.3
	GPQA Diamond (Pass@1)	65.0	49.9	59.1	60.0	75.7	71.5
	SimpleQA (Correct)	28.4	38.2	24.9	7.0	47.0	30.1
	FRAMES (Acc.)	72.5	80.5	73.3	76.9	-	82.5
	AlpacaEval2.0 (LC-winrate)	52.0	51.1	70.0	57.8	-	87.6
	ArenaHard (GPT-4-1106)	85.2	80.4	85.5	92.0	-	92.3
Code	LiveCodeBench (Pass@1-COT)	38.9	32.9	36.2	53.8	63.4	65.9
	Codeforces (Percentile)	20.3	23.6	58.7	93.4	96.6	96.3
	Codeforces (Rating)	717	759	1134	1820	2061	2029
	SWE Verified (Resolved)	50.8	38.8	42.0	41.6	48.9	49.2
	Aider-Polyglot (Acc.)	45.3	16.0	49.6	32.9	61.7	53.3
Math	AIME 2024 (Pass@1)	16.0	9.3	39.2	63.6	79.2	79.8
	MATH-500 (Pass@1)	78.3	74.6	90.2	90.0	96.4	97.3
	CNMO 2024 (Pass@1)	13.1	10.8	43.2	67.6	-	78.8
Chinese	CLUEWSC (EM)	85.4	87.9	90.9	89.9	-	92.8
	C-Eval (EM)	76.7	76.0	86.5	68.9	-	91.8
	C-SimpleQA (Correct)	55.4	58.7	68.0	40.3	-	63.7

Table 4 | DeepSeek-R1 v1.7, 72B, 1/2, 1/2, 1/2

/ i 2 k%o h RL > Z j 1/2 j 2 j 2 1/2 j 1/2 1/2 %o h RL DeepSeek-R1 i j 2 1/2... i 70%,
i n 2

DeepSeek-R1 (IF-Eval) 在 10 个任务上取得 85.2% 的平均得分，优于 GPT-4o 的 82.1%。在 SFT 和 RL 数据集上，DeepSeek-R1 的得分分别为 88.5% 和 89.1%。在 AlpacaEval2.0 和 ArenaHard 数据集上，DeepSeek-R1 的得分分别为 88.5% 和 89.1%。在 10 个任务上，DeepSeek-R1 的得分分别为 85.2%、86.1%、87.0%、87.9%、88.8%、89.7%、90.6%、91.5%、92.4%、93.3%。在 10 个任务上，DeepSeek-R1 的得分分别为 85.2%、86.1%、87.0%、87.9%、88.8%、89.7%、90.6%、91.5%、92.4%、93.3%。

(p f i ě ½ DeepSeek-R1„ h° OpenAI-o1-1217i § ½' E...Š v i ě ½ (—i ě ½½ , LiveCodeBenchECodeforces i Ž ½ „ ! < ; i ě ½½½ ě ½½ ě ½½ ě ½½ i ě ½ i ě ½ OpenAI-o1-1217(Alder Ž DeepSeek-R1 F (SWE Verified h° i § ½ i ě ½½• @i § R- i p ě ½½½ž DeepSeek-R1„ i ě ½' i ě ½ * H, - —O9„ i M ě ½½ ě ½½ P

Evaluation result Prompts Following the setup in DeepSeek-V3, standard benchmarks such as MMLU, DROP, GPQA Diamond, and SimpleQA are evaluated using prompts from the simple-evals framework. For MMLU-Redux, we adopt the Zero-Eval prompt format (Lin, 2024) in a zero-shot setting. In terms of MMLU-Pro, C-Eval and CLUE-WSC, since the original prompts are few-shot, we slightly modify the prompt to the zero-shot setting. The CoT in few-shot may hurt the performance of DeepSeek-R1. Other datasets follow their original evaluation protocols with default prompts provided by their creators. For code and math benchmarks, the HumanEval-Mul dataset covers eight mainstream programming languages (Python, Java, C++, C#, JavaScript, TypeScript, PHP, and Bash). Model performance on LiveCodeBench is evaluated using CoT format, with data collected between August 2024 and January 2025. The Codeforces dataset is evaluated using problems from 10 Div.2 contests along with expert-crafted test cases, after which the expected ratings and percentages of competitors are calculated. SWE-Bench verified results are obtained via the agentless framework (Xia et al., 2024). AIDER-related benchmarks are measured using a "diff" format. DeepSeek-R1 outputs are capped at a maximum of 32,768 tokens for each benchmark.

Baselines We conduct comprehensive evaluations against several strong baselines, including DeepSeek-V3, Claude-Sonnet-3.5-1022, GPT-4o-0513, OpenAI-o1-mini, and OpenAI-o1-1217. Since accessing the OpenAI-o1-1217 API is challenging in mainland China, we report its performance based on official reports. For distilled models, we also compare the open-source model QwQ-32B-Preview (Qwen, 2024a).

15

Model	AIME 2024		MATH-500	GPOA Diamond	LiveCode Bench	CodeForces
	pass@1	cons@64	pass@1	pass@1	pass@1	rating
GPT-4o-0513	9.3	13.4	74.6	49.9	32.9	759
Claude-3.5-Sonnet-1022	16.0	26.7	78.3	65.0	38.9	717
OpenAI-o1-mini	63.6	80.0	90.0	60.0	53.8	1820
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9	1316
DeepSeek-R1-Distill-Qwen-1.5B	28.9	52.7	83.9	33.8	16.9	954
DeepSeek-R1-Distill-Qwen-7B	55.5	83.3	92.8	49.1	37.6	1189
DeepSeek-R1-Distill-Qwen-14B	69.7	80.0	93.9	59.1	53.1	1481
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2	1691
DeepSeek-R1-Distill-Llama-8B	50.4	80.0	89.1	49.0	39.6	1205
DeepSeek-R1-Distill-Llama-70B	70.0	86.7	94.5	65.2	57.5	1633

, h5@: € U, • DeepSeek-R1, " i ½ ½ H, DeepSeek-R1-7B s DeepSeek-R1-Distill-Qwen-7B i ½ ½ < € i ½ ½ @ i ½ ½ ½ ½ ½ ½ ! < , GPT-4o-0513 DeepSeek-R1-14B(@ i ½ ½ ... i ½ ½ QwQ-32B-Preview DeepSeek-R1-32B(DeepSeek-R1-70B(' p i ½ ½ ½ ½ W... i ½ ½ ½ mini i ½ ½ geU: †, •, „: ' \> d i ½ ½ ½: f`" (Z i ½ ½ • ! < i ½ ½ ½ > W, i ½ ½ ½ G i ½ ½ ½ ½ i ½ ½ ½ ½ ½ ½ " i ½ ½ ½ ½ ½ ½ ½ ½ € U, SFT, • ! < „ i ½ ½ ½

4.1. Distillation v.s. Reinforcement Learning

4.1. Distillation v.s. Reinforcement Learning

Table 6 j i 2 1/2 < E: f` ! < (" i 8 V2j 1/2 1/2 j 1/2

: + ĩ 72ĩ 8 72ĭ 72 ĩ ħ Qwen-32B-Base • (p f ĩ ĩ ĩ STEM p n ĩ ĩ 72' ĩ 72
„ : f` - ĩ ĩ 72ĩ 72ĩ 72ĩ 72ĩ ĩ ĩ ĩ DeepSeek-R1-Zero-Qwen-32B h 6 - „ ž
Eĩ e 72: ĩ ĩ 727272: f` - ĩ ĩ 72 32Bĩ @ 72< „ ' ĩ ĩ 72wQ-32B-Preview ĩ S 726
ĩ ĩ DeepSeek-R1 „ —O„ DeepSeek-R1-Distill-Qwen-32B (@ ĩ ĩ 7272 ĩ 72h° >
Ž DeepSeek-R1-Zero-Qwen-32B

is then calculated as

pass@1 = \frac{1}{9} \sum_{l=1}^9 \mathbb{1}_{r_l > \tau}

where τ_l denotes the correctness of the l -th response. This method provides more reliable performance estimates. For AIME 2024, we also report consensus (majority vote) results (Wang et al., 2022) using 64 samples, denoted as cons@64.

3.1. DeepSeek-R1 Evaluation

Benchmark (Metric)	Claude-3.5-Sonnet-1022	GPT-4o-0513	DeepSeek-V3	OpenAI-o1-mini	OpenAI-o1-1217	DeepSeek-R1
Architecture	-	-	MoE	-	-	MoE
# Activated Params	-	-	37B	-	-	37B
# Total Params	-	-	671B	-	-	671B
English	MMLU (Pass@1)	88.3	87.2	88.5	85.2	91.8
	MMLU-Redux (EM)	88.9	88.0	89.1	86.7	-
	MMLU-Pro (EM)	78.0	72.6	75.9	80.3	-
	DROP (3-shot F1)	88.3	83.7	91.6	83.9	90.2
	IF-Eval (Prompt Strict)	86.5	84.3	86.1	84.8	-
	GPOA Diamond (Pass@1)	65.0	49.9	59.1	60.0	75.7
	SimpleQA (Correct)	28.4	38.2	24.9	7.0	47.0
	FRAMES (Acc.)	72.5	80.5	73.3	76.9	-
	AlpacaEval2.0 (LC-winrate)	52.0	51.1	70.0	57.8	-
Code	ArenaHard (GPT-4-1106)	85.2	80.4	85.5	92.0	-
	LiveCodeBench (Pass@1-COT)	38.9	32.9	36.2	53.8	63.4
	Codeforces (Percentile)	20.3	23.6	58.7	93.4	96.6
	Codeforces (Rating)	717	759	1134	1820	2061
	SWE Verified (Resolved)	50.8	38.8	42.0	41.6	48.9
	Aider-Polyglot (Acc.)	45.3	16.0	49.6	32.9	61.7
Math	AIME 2024 (Pass@1)	16.0	9.3	39.2	63.6	79.2
	MATH-500 (Pass@1)	78.3	74.6	90.2	90.0	96.4
	CNMO 2024 (Pass@1)	13.1	10.8	43.2	67.6	-
Chinese	CLUEWSC (EM)	85.4	87.9	90.9	89.9	-
	C-Eval (EM)	76.7	76.0	86.5	68.9	-
	C-SimpleQA (Correct)	55.4	58.7	68.0	40.3	-

Table 4 j Comparison between DeepSeek-R1 and other representative models.

For education-oriented knowledge benchmarks such as MMLU, MMLU-Pro, and GPOA Diamond, DeepSeek-R1 demonstrates superior performance compared to DeepSeek-V3. This improvement is primarily attributed to enhanced accuracy in STEM-related questions, where significant gains are achieved through large-scale reinforcement learning. Additionally, DeepSeek-R1 excels on FRAMES, a long-context-dependent QA task, showcasing its strong document analysis capabilities. This highlights the potential of reasoning models in AI-driven search and data analysis tasks. On the factual benchmark SimpleQA, DeepSeek-R1 outperforms DeepSeek-V3, demonstrating its capability in handling fact-based queries. A similar trend is observed where OpenAI-o1 surpasses GPT-4o on this benchmark. However, DeepSeek-R1 performs worse than

OpenAI-o1 on the MMLU-Pro benchmark. In the domain of code generation, DeepSeek-R1 shows competitive performance with OpenAI-o1, though it is slightly outperformed by GPT-4o on the LiveCodeBench and SWE Verified benchmarks. Overall, DeepSeek-R1 demonstrates strong performance across a wide range of benchmarks, particularly in the areas of reasoning and document analysis.

4.2. Unsuccessful Attempts

During the evaluation process, we encountered several challenges that prevented us from achieving the desired results. One major issue was the limited availability of certain benchmarks, which restricted our ability to fully assess the model's performance in specific domains. Additionally, the high variability in the results across different runs and samples made it difficult to draw consistent conclusions. We also faced challenges related to the computational resources required for the evaluation, which sometimes led to incomplete results or timeouts.

One of the primary reasons for the unsuccessful attempts was the lack of sufficient data for certain benchmarks, which prevented us from training the model effectively. Another significant challenge was the complexity of the tasks, particularly in the domain of code generation, where the model often struggled to produce accurate and functional code. We also encountered issues related to the evaluation process itself, such as the inconsistency in the results and the difficulty in interpreting the model's outputs. Despite these challenges, we were able to identify several areas for improvement and implement changes to enhance the model's performance in the future.

One of the key factors contributing to the unsuccessful attempts was the limited availability of certain benchmarks, which restricted our ability to fully assess the model's performance in specific domains. Additionally, the high variability in the results across different runs and samples made it difficult to draw consistent conclusions. We also faced challenges related to the computational resources required for the evaluation, which sometimes led to incomplete results or timeouts. Despite these challenges, we were able to identify several areas for improvement and implement changes to enhance the model's performance in the future.

One of the primary reasons for the unsuccessful attempts was the lack of sufficient data for certain benchmarks, which prevented us from training the model effectively. Another significant challenge was the complexity of the tasks, particularly in the domain of code generation, where the model often struggled to produce accurate and functional code. We also encountered issues related to the evaluation process itself, such as the inconsistency in the results and the difficulty in interpreting the model's outputs. Despite these challenges, we were able to identify several areas for improvement and implement changes to enhance the model's performance in the future.

One of the key factors contributing to the unsuccessful attempts was the limited availability of certain benchmarks, which restricted our ability to fully assess the model's performance in specific domains. Additionally, the high variability in the results across different runs and samples made it difficult to draw consistent conclusions. We also faced challenges related to the computational resources required for the evaluation, which sometimes led to incomplete results or timeouts. Despite these challenges, we were able to identify several areas for improvement and implement changes to enhance the model's performance in the future.

DeepSeek-V3 on the Chinese SimpleQA benchmark, primarily due to its tendency to refuse answering certain queries after safety RL. Without safety RL, DeepSeek-R1 could achieve an accuracy of over 70%.

DeepSeek-R1 also delivers impressive results on IF-Eval, a benchmark designed to assess a model’s ability to follow format instructions. These improvements can be linked to the inclusion of instruction-following data during the final stages of supervised fine-tuning (SFT) and RL training. Furthermore, remarkable performance is observed on AlpacaEval2.0 and ArenaHard, indicating DeepSeek-R1’s strengths in writing tasks and open-domain question answering. Its significant outperformance of DeepSeek-V3 underscores the generalization benefits of large-scale RL, which not only boosts reasoning capabilities but also improves performance across diverse domains. Moreover, the summary lengths generated by DeepSeek-R1 are concise, with an average of 689 tokens on ArenaHard and 2,218 characters on AlpacaEval 2.0. This indicates that DeepSeek-R1 avoids introducing length bias during GPT-based evaluations, further solidifying its robustness across multiple tasks.

On math tasks, DeepSeek-R1 demonstrates performance on par with OpenAI-o1-1217, surpassing other models by a large margin. A similar trend is observed on coding algorithm tasks, such as LiveCodeBench and Codeforces, where reasoning-focused models dominate these benchmarks. On engineering-oriented coding tasks, OpenAI-o1-1217 outperforms DeepSeek-R1 on Aider but achieves comparable performance on SWE Verified. We believe the engineering performance of DeepSeek-R1 will improve in the next version, as the amount of related RL training data currently remains very limited.

3.2. Distilled Model Evaluation

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCode Bench	CodeForces
	pass@1	cons@64	pass@1	pass@1	pass@1	rating
GPT-4o-0513	9.3	13.4	74.6	49.9	32.9	759
Claude-3.5-Sonnet-1022	16.0	26.7	78.3	65.0	38.9	717
OpenAI-o1-mini	63.6	80.0	90.0	60.0	53.8	1820
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9	1316
DeepSeek-R1-Distill-Qwen-1.5B	28.9	52.7	83.9	33.8	16.9	954
DeepSeek-R1-Distill-Qwen-7B	55.5	83.3	92.8	49.1	37.6	1189
DeepSeek-R1-Distill-Qwen-14B	69.7	80.0	93.9	59.1	53.1	1481
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2	1691
DeepSeek-R1-Distill-Llama-8B	50.4	80.0	89.1	49.0	39.6	1205
DeepSeek-R1-Distill-Llama-70B	70.0	86.7	94.5	65.2	57.5	1633

Table 5 j Comparison of DeepSeek-R1 distilled models and other comparable models on reasoning-related benchmarks.

As shown in Table 5, simply distilling DeepSeek-R1’s outputs enables the efficient DeepSeek-

5. Conclusion, Limitations, and Future Work

(, i x- i z t i ½ f` ž : ! < " i z t 44>(449 '1 /4@= i b ½
i z ½ RLi z ½ • VŽ . / " p n (i z ½ h° i z ½ 44>(449 '1 i z ½'
) (. / " p n v i z ½ R½® i z ½ i z ½ ž ° + OpenAI-o1-1217 i z ½'
i z ½

[illegible]
$$(\ast e \quad i \frac{1}{2} (i \frac{1}{2} i \frac{1}{2} i \frac{1}{2}) \frac{1}{2}) (449 \quad 1 i \frac{1}{2} v \cdot D$$

- [illegible]

References

AI@Meta. Llama 3.1 model card, 2024. URL https://github.com/meta-llama/llama-models/blob/main/models/llama3.1/MODEL_CARD.md.

Anthropic. Claude 3.5 sonnet, 2024. URL <https://www.anthropic.com/news/claude-3-5-sonnet>.

M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A.

R1-7B (i.e., DeepSeek-R1-Distill-Qwen-7B, abbreviated similarly below) to outperform non-reasoning models like GPT-4o-0513 across the board. DeepSeek-R1-14B surpasses QwQ-32B-Preview on all evaluation metrics, while DeepSeek-R1-32B and DeepSeek-R1-70B significantly exceed o1-mini on most benchmarks. These results demonstrate the strong potential of distillation. Additionally, we found that applying RL to these distilled models yields significant further gains. We believe this warrants further exploration and therefore present only the results of the simple SFT-distilled models here.

4. Discussion

4.1. Distillation v.s. Reinforcement Learning

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCodeBench
	pass@1	cons@64	pass@1	pass@1	pass@1
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9
DeepSeek-R1-Zero-Qwen-32B	47.0	60.0	91.6	55.0	40.2
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2

Table 6 j Comparison of distilled and RL Models on Reasoning-Related Benchmarks.

In Section 3.2, we can see that by distilling DeepSeek-R1, the small model can achieve impressive results. However, there is still one question left: can the model achieve comparable performance through the large-scale RL training discussed in the paper without distillation?

To answer this question, we conduct large-scale RL training on Qwen-32B-Base using math, code, and STEM data, training for over 10K steps, resulting in DeepSeek-R1-Zero-Qwen-32B. The experimental results, shown in Table 6, demonstrate that the 32B base model, after large-scale RL training, achieves performance on par with QwQ-32B-Preview. However, DeepSeek-R1-Distill-Qwen-32B, which is distilled from DeepSeek-R1, performs significantly better than DeepSeek-R1-Zero-Qwen-32B across all benchmarks.

Therefore, we can draw two conclusions: First, distilling more powerful models into smaller ones yields excellent results, whereas smaller models relying on the large-scale RL mentioned in this paper require enormous computational power and may not even achieve the performance of distillation. Second, while distillation strategies are both economical and effective, advancing beyond the boundaries of intelligence may still require more powerful base models and larger-scale reinforcement learning.

4.2. Unsuccessful Attempts

In the early stages of developing DeepSeek-R1, we also encountered failures and setbacks along the way. We share our failure experiences here to provide insights, but this does not imply that

Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021. URL <https://arxiv.org/abs/2107.03374>.

A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Y. Dubois, B. Galambosi, P. Liang, and T. B. Hashimoto. Length-controlled alpacaEval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.

X. Feng, Z. Wan, M. Wen, S. M. McAleer, Y. Wen, W. Zhang, and J. Wang. Alphazero-like tree-search can guide large language model decoding and training, 2024. URL <https://arxiv.org/abs/2309.17179>.

L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization, 2022. URL <https://arxiv.org/abs/2210.10760>.

A. P. Gema, J. O. J. Leang, G. Hong, A. Devoto, A. C. M. Mancino, R. Saxena, X. He, Y. Zhao, X. Du, M. R. G. Madani, C. Barale, R. McHardy, J. Harris, J. Kaddour, E. van Krieken, and P. Minervini. Are we done with mmlu? *CoRR*, abs/2406.04127, 2024. URL <https://doi.org/10.48550/arXiv.2406.04127>.

Google. Our next-generation model: Gemini 1.5, 2024. URL <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024>.

Y. He, S. Li, J. Liu, Y. Tan, W. Wang, H. Huang, X. Bu, H. Guo, C. Hu, B. Zheng, et al. Chinese simpleqa: A chinese factuality evaluation for large language models. *arXiv preprint arXiv:2411.07140*, 2024.

D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

Y. Huang, Y. Bai, Z. Zhu, J. Zhang, J. Zhang, T. Su, J. Liu, C. Lv, Y. Zhang, J. Lei, et al. C-Eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *arXiv preprint arXiv:2305.08322*, 2023.

N. Jain, K. Han, A. Gu, W. Li, F. Yan, T. Zhang, S. Wang, A. Solar-Lezama, K. Sen, and I. Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *CoRR*, abs/2403.07974, 2024. URL <https://doi.org/10.48550/arXiv.2403.07974>.

S. Krishna, K. Krishna, A. Mohananey, S. Schwarcz, A. Stambler, S. Upadhyay, and M. Faruqui. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. *CoRR*,

these approaches are incapable of developing effective reasoning models.

Process Reward Model (PRM) PRM is a reasonable method to guide the model toward better approaches for solving reasoning tasks (Lightman et al., 2023; Uesato et al., 2022; Wang et al., 2023). However, in practice, PRM has three main limitations that may hinder its ultimate success. First, it is challenging to explicitly define a fine-grain step in general reasoning. Second, determining whether the current intermediate step is correct is a challenging task. Automated annotation using models may not yield satisfactory results, while manual annotation is not conducive to scaling up. Third, once a model-based PRM is introduced, it inevitably leads to reward hacking (Gao et al., 2022), and retraining the reward model needs additional training resources and it complicates the whole training pipeline. In conclusion, while PRM demonstrates a good ability to rerank the top-N responses generated by the model or assist in guided search (Snell et al., 2024), its advantages are limited compared to the additional computational overhead it introduces during the large-scale reinforcement learning process in our experiments.

Monte Carlo Tree Search (MCTS) Inspired by AlphaGo (Silver et al., 2017b) and AlphaZero (Silver et al., 2017a), we explored using Monte Carlo Tree Search (MCTS) to enhance test-time compute scalability. This approach involves breaking answers into smaller parts to allow the model to explore the solution space systematically. To facilitate this, we prompt the model to generate multiple tags that correspond to specific reasoning steps necessary for the search. For training, we first use collected prompts to find answers via MCTS guided by a pre-trained value model. Subsequently, we use the resulting question-answer pairs to train both the actor model and the value model, iteratively refining the process.

However, this approach encounters several challenges when scaling up the training. First, unlike chess, where the search space is relatively well-defined, token generation presents an exponentially larger search space. To address this, we set a maximum extension limit for each node, but this can lead to the model getting stuck in local optima. Second, the value model directly influences the quality of generation since it guides each step of the search process. Training a fine-grained value model is inherently difficult, which makes it challenging for the model to iteratively improve. While AlphaGo’s core success relied on training a value model to progressively enhance its performance, this principle proves difficult to replicate in our setup due to the complexities of token generation.

In conclusion, while MCTS can improve performance during inference when paired with a pre-trained value model, iteratively boosting model performance through self-search remains a significant challenge.

abs/2409.12941, 2024. doi: 10.48550/ARXIV.2409.12941. URL <https://doi.org/10.48550/arXiv.2409.12941>.

A. Kumar, V. Zhuang, R. Agarwal, Y. Su, J. D. Co-Reyes, A. Singh, K. Baumli, S. Iqbal, C. Bishop, R. Roelofs, et al. Training language models to self-correct via reinforcement learning. [arXiv preprint arXiv:2409.12917](https://arxiv.org/abs/2409.12917), 2024.

H. Li, Y. Zhang, F. Koto, Y. Yang, H. Zhao, Y. Gong, N. Duan, and T. Baldwin. CMMLU: Measuring massive multitask language understanding in Chinese. [arXiv preprint arXiv:2306.09212](https://arxiv.org/abs/2306.09212), 2023.

T. Li, W.-L. Chiang, E. Frick, L. Dunlap, T. Wu, B. Zhu, J. E. Gonzalez, and I. Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. [arXiv preprint arXiv:2406.11939](https://arxiv.org/abs/2406.11939), 2024.

H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe. Let’s verify step by step. [arXiv preprint arXiv:2305.20050](https://arxiv.org/abs/2305.20050), 2023.

B. Y. Lin. ZeroEval: A Unified Framework for Evaluating Language Models, July 2024. URL <https://github.com/WildEval/ZeroEval>.

MAA. American invitational mathematics examination - aime. In [American Invitational Mathematics Examination - AIME 2024](https://maa.org/math-competitions/american-invitational-mathematics-examination-aime), February 2024. URL <https://maa.org/math-competitions/american-invitational-mathematics-examination-aime>.

OpenAI. Hello GPT-4o, 2024a. URL <https://openai.com/index/hello-gpt-4o/>.

OpenAI. Learning to reason with llms, 2024b. URL <https://openai.com/index/learning-to-reason-with-llms/>.

OpenAI. Introducing SimpleQA, 2024c. URL <https://openai.com/index/introducing-simpleqa/>.

OpenAI. Introducing SWE-bench verified we’re releasing a human-validated subset of swe-bench that more, 2024d. URL <https://openai.com/index/introducing-swe-bench-verified/>.

Qwen. Qwq: Reflect deeply on the boundaries of the unknown, 2024a. URL <https://qwenlm.github.io/blog/qwq-32b-preview/>.

Qwen. Qwen2.5: A party of foundation models, 2024b. URL <https://qwenlm.github.io/blog/qwen2.5/>.

D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. [arXiv preprint arXiv:2311.12022](https://arxiv.org/abs/2311.12022), 2023.

5. Conclusion, Limitations, and Future Work

In this work, we share our journey in enhancing model reasoning abilities through reinforcement learning. DeepSeek-R1-Zero represents a pure RL approach without relying on cold-start data, achieving strong performance across various tasks. DeepSeek-R1 is more powerful, leveraging cold-start data alongside iterative RL fine-tuning. Ultimately, DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on a range of tasks.

We further explore distillation the reasoning capability to small dense models. We use DeepSeek-R1 as the teacher model to generate 800K training samples, and fine-tune several small dense models. The results are promising: DeepSeek-R1-Distill-Qwen-1.5B outperforms GPT-4o and Claude-3.5-Sonnet on math benchmarks with 28.9% on AIME and 83.9% on MATH. Other dense models also achieve impressive results, significantly outperforming other instruction-tuned models based on the same underlying checkpoints.

In the future, we plan to invest in research across the following directions for DeepSeek-R1.

- **General Capability:** Currently, the capabilities of DeepSeek-R1 fall short of DeepSeek-V3 in tasks such as function calling, multi-turn, complex role-playing, and JSON output. Moving forward, we plan to explore how long CoT can be leveraged to enhance tasks in these fields.
- **Language Mixing:** DeepSeek-R1 is currently optimized for Chinese and English, which may result in language mixing issues when handling queries in other languages. For instance, DeepSeek-R1 might use English for reasoning and responses, even if the query is in a language other than English or Chinese. We aim to address this limitation in future updates.
- **Prompting Engineering:** When evaluating DeepSeek-R1, we observe that it is sensitive to prompts. Few-shot prompting consistently degrades its performance. Therefore, we recommend users directly describe the problem and specify the output format using a zero-shot setting for optimal results.
- **Software Engineering Tasks:** Due to the long evaluation times, which impact the efficiency of the RL process, large-scale RL has not been applied extensively in software engineering tasks. As a result, DeepSeek-R1 has not demonstrated a huge improvement over DeepSeek-V3 on software engineering benchmarks. Future versions will address this by implementing rejection sampling on software engineering data or incorporating asynchronous evaluations during the RL process to improve efficiency.

- Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, M. Zhang, Y. Li, Y. Wu, and D. Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. [arXiv preprint arXiv:2402.03300](#), 2024.
- D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. P. Lillicrap, K. Simonyan, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. [CoRR](#), abs/1712.01815, 2017a. URL <http://arxiv.org/abs/1712.01815>.
- D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. P. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis. Mastering the game of go without human knowledge. [Nat.](#), 550(7676):354–359, 2017b. doi: 10.1038/NATURE24270. URL <https://doi.org/10.1038/nature24270>.
- C. Snell, J. Lee, K. Xu, and A. Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.
- T. Trinh, Y. Wu, Q. Le, H. He, and T. Luong. Solving olympiad geometry without human demonstrations. [Nature](#), 2024. doi: 10.1038/s41586-023-06747-5.
- J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins. Solving math word problems with process-and outcome-based feedback. [arXiv preprint arXiv:2211.14275](#), 2022.
- P. Wang, L. Li, Z. Shao, R. Xu, D. Dai, Y. Li, D. Chen, Y. Wu, and Z. Sui. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. [arXiv preprint arXiv:2312.08935](#), 2023.
- X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. [arXiv preprint arXiv:2203.11171](#), 2022.
- Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. [CoRR](#), abs/2406.01574, 2024. URL <https://doi.org/10.48550/arXiv.2406.01574>.
- C. S. Xia, Y. Deng, S. Dunn, and L. Zhang. Agentless: Demystifying llm-based software engineering agents. [arXiv preprint](#), 2024.
- H. Xin, Z. Z. Ren, J. Song, Z. Shao, W. Zhao, H. Wang, B. Liu, L. Zhang, X. Lu, Q. Du, W. Gao, Q. Zhu, D. Yang, Z. Gou, Z. F. Wu, F. Luo, and C. Ruan. Deepseek-prover-v1.5: Harnessing

References

- AI@Meta. Llama 3.1 model card, 2024. URL https://github.com/meta-llama/llama-models/blob/main/models/llama3_1/MODEL_CARD.md.
- Anthropic. Claude 3.5 sonnet, 2024. URL <https://www.anthropic.com/news/claude-3-5-sonnet>.
- M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021. URL <https://arxiv.org/abs/2107.03374>.
- A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Y. Dubois, B. Galambosi, P. Liang, and T. B. Hashimoto. Length-controlled alpacaEval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- X. Feng, Z. Wan, M. Wen, S. M. McAleer, Y. Wen, W. Zhang, and J. Wang. Alphazero-like tree-search can guide large language model decoding and training, 2024. URL <https://arxiv.org/abs/2309.17179>.
- L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization, 2022. URL <https://arxiv.org/abs/2210.10760>.
- A. P. Gema, J. O. J. Leang, G. Hong, A. Devoto, A. C. M. Mancino, R. Saxena, X. He, Y. Zhao, X. Du, M. R. G. Madani, C. Barale, R. McHardy, J. Harris, J. Kaddour, E. van Krieken, and P. Minervini. Are we done with mmlu? *CoRR*, abs/2406.04127, 2024. URL <https://doi.org/10.48550/arXiv.2406.04127>.
- Google. Our next-generation model: Gemini 1.5, 2024. URL <https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024>.
- Y. He, S. Li, J. Liu, Y. Tan, W. Wang, H. Huang, X. Bu, H. Guo, C. Hu, B. Zheng, et al. Chinese simpleQA: A chinese factuality evaluation for large language models. *arXiv preprint arXiv:2411.07140*, 2024.

proof assistant feedback for reinforcement learning and monte-carlo tree search, 2024. URL <https://arxiv.org/abs/2408.08152>.

J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. [arXiv preprint arXiv:2009.03300](#), 2020.

Y. Huang, Y. Bai, Z. Zhu, J. Zhang, J. Zhang, T. Su, J. Liu, C. Lv, Y. Zhang, J. Lei, et al. C-Eval: A multi-level multi-discipline chinese evaluation suite for foundation models. [arXiv preprint arXiv:2305.08322](#), 2023.

N. Jain, K. Han, A. Gu, W. Li, F. Yan, T. Zhang, S. Wang, A. Solar-Lezama, K. Sen, and I. Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. [CoRR](#), abs/2403.07974, 2024. URL <https://doi.org/10.48550/arXiv.2403.07974>.

S. Krishna, K. Krishna, A. Mohananey, S. Schwarcz, A. Stambler, S. Upadhyay, and M. Faruqui. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. [CoRR](#), abs/2409.12941, 2024. doi: 10.48550/ARXIV.2409.12941. URL <https://doi.org/10.48550/arXiv.2409.12941>.

A. Kumar, V. Zhuang, R. Agarwal, Y. Su, J. D. Co-Reyes, A. Singh, K. Baumli, S. Iqbal, C. Bishop, R. Roelofs, et al. Training language models to self-correct via reinforcement learning. [arXiv preprint arXiv:2409.12917](#), 2024.

H. Li, Y. Zhang, F. Koto, Y. Yang, H. Zhao, Y. Gong, N. Duan, and T. Baldwin. CMMLU: Measuring massive multitask language understanding in Chinese. [arXiv preprint arXiv:2306.09212](#), 2023.

T. Li, W.-L. Chiang, E. Frick, L. Dunlap, T. Wu, B. Zhu, J. E. Gonzalez, and I. Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. [arXiv preprint arXiv:2406.11939](#), 2024.

H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe. Let’s verify step by step. [arXiv preprint arXiv:2305.20050](#), 2023.

B. Y. Lin. ZeroEval: A Unified Framework for Evaluating Language Models, July 2024. URL <https://github.com/WildEval/ZeroEval>.

MAA. American invitational mathematics examination - aime. In [American Invitational Mathematics Examination - AIME 2024](#), February 2024. URL <https://maa.org/math-competitions/american-invitational-mathematics-examination-aime>.

OpenAI. Hello GPT-4o, 2024a. URL <https://openai.com/index/hello-gpt-4o/>.

OpenAI. Learning to reason with llms, 2024b. URL <https://openai.com/index/learning-to-reason-with-llms/>.

OpenAI. Introducing SimpleQA, 2024c. URL <https://openai.com/index/introducing-simpleqa/>.

Appendix

A. Contributions and Acknowledgments

8½.	—i ½
i ½ ½	4• j
h• i ½	Wi ½
i	i ½
< ½	H ÷
i ½	Ni ½
• i ½	H.
1 / j	• l
l c	* %
< i ½	½ \
½	Hg
.	N%
Y t i ½	i ½ ½
4±	N i ½
4x	Ho
i ½	• o 3
i ½	, ½
NS	± ½
½ ½	N>™
	! J.L.
!	* *
1 #	• e
i ½	Hi ½
< i ½	c i ½
4Z...	i ½ ½
-	\$ i ½
b i ½	½ ½
u	s•
" *	½ # ½
. i ½	Y ½
4' i	< ½
H• i ½	Pj
c p	u®
N *	< i ½

OpenAI. Introducing SWE-bench verified we re releasing a human-validated subset of swe-bench that more, 2024d. URL <https://openai.com/index/introducing-swe-bench-verified/>.

Qwen. Qwq: Reflect deeply on the boundaries of the unknown, 2024a. URL <https://qwenlm.github.io/blog/qwq-32b-preview/>.

Qwen. Qwen2.5: A party of foundation models, 2024b. URL <https://qwenlm.github.io/blog/qwen2.5>.

D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. arXiv preprint arXiv:2311.12022, 2023.

Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, M. Zhang, Y. Li, Y. Wu, and D. Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300, 2024.

D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. P. Lillicrap, K. Simonyan, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. CoRR, abs/1712.01815, 2017a. URL <http://arxiv.org/abs/1712.01815>.

D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. P. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis. Mastering the game of go without human knowledge. Nat., 550(7676):354–359, 2017b. doi: 10.1038/NATURE24270. URL <https://doi.org/10.1038/nature24270>.

C. Snell, J. Lee, K. Xu, and A. Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.

T. Trinh, Y. Wu, Q. Le, H. He, and T. Luong. Solving olympiad geometry without human demonstrations. Nature, 2024. doi: 10.1038/s41586-023-06747-5.

J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins. Solving math word problems with process-and outcome-based feedback. arXiv preprint arXiv:2211.14275, 2022.

P. Wang, L. Li, Z. Shao, R. Xu, D. Dai, Y. Li, D. Chen, Y. Wu, and Z. Sui. Math-shepherd: A label-free step-by-step verifier for llms in mathematical reasoning. arXiv preprint arXiv:2312.08935, 2023.

X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171, 2022.

=

• $\frac{1}{2}$

$\frac{1}{2}\frac{1}{2}\frac{1}{2}$

$\frac{1}{2}$

N

%o

h $\frac{1}{2}$

N

< x >

N

O•

$\frac{1}{2}\frac{1}{2}\frac{1}{2}\frac{1}{2}\frac{1}{2}$

O

< $\sim \frac{1}{2}$

H• $\frac{1}{2}$

$\setminus \frac{1}{2}\frac{1}{2}\frac{1}{2}$

[^ * *

^ ~

X• $\frac{1}{2}$

< $\frac{1}{2}$

R.J. H

R.L. $\frac{1}{2}$

H,

b i

h $\frac{1}{2}$

Hq $\frac{1}{2}$

$\frac{1}{2}\frac{1}{2}\frac{1}{2}\frac{1}{2}$

< $\frac{1}{2}$

Y 4

h z $\frac{1}{2}$

X $\frac{1}{2}$

S.S. N

h =

4 R

$\frac{1}{2}\frac{1}{2}\frac{1}{2}\frac{1}{2}$

' >

$\frac{1}{2}$

Y) $\frac{1}{2}$

T. <

$\frac{1}{2}$

$\frac{1}{2}$

• $\frac{1}{2}$

$\frac{1}{2}$

Y $\frac{1}{2}$

Wentao

W.L. —

%o

S

< S $\frac{1}{2}$

H .

B >

#

#

" #

t ...

h $\frac{1}{2}$

N $\frac{1}{2}$

$\frac{1}{2}$

$\frac{1}{2}$

X.Q. N

$\frac{1}{2}\frac{1}{2}\frac{1}{2}\frac{1}{2}$

^ S ~

HS $\frac{1}{2}$

YS $\frac{1}{2}$

< $\frac{1}{2}$

< $\frac{1}{2}$

h # $\frac{1}{2}$

< HV

UK

Y.K. N

Y.Q. <

Y.X. O

3

• $\frac{1}{2}$

Nv

uv

Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *CoRR*, abs/2406.01574, 2024. URL <https://doi.org/10.48550/arXiv.2406.01574>.

C. S. Xia, Y. Deng, S. Dunn, and L. Zhang. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint*, 2024.

H. Xin, Z. Z. Ren, J. Song, Z. Shao, W. Zhao, H. Wang, B. Liu, L. Zhang, X. Lu, Q. Du, W. Gao, Q. Zhu, D. Yang, Z. Gou, Z. F. Wu, F. Luo, and C. Ruan. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search, 2024. URL <https://arxiv.org/abs/2408.08152>.

J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

Y. Ni, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *CoRR*, abs/2406.01574, 2024. URL <https://doi.org/10.48550/arXiv.2406.01574>.

C. S. Xia, Y. Deng, S. Dunn, and L. Zhang. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint*, 2024.

H. Xin, Z. Z. Ren, J. Song, Z. Shao, W. Zhao, H. Wang, B. Liu, L. Zhang, X. Lu, Q. Du, W. Gao, Q. Zhu, D. Yang, Z. Gou, Z. F. Wu, F. Luo, and C. Ruan. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search, 2024. URL <https://arxiv.org/abs/2408.08152>.

J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

Appendix

A. Contributions and Acknowledgments

Core Contributors	Fangyun Lin
Daya Guo	Fucong Dai
Dejian Yang	Fuli Luo*
Haowei Zhang	Guangbo Hao
Junxiao Song	Guanting Chen
Ruoyu Zhang	Guowei Li
Runxin Xu	H. Zhang
Qihao Zhu	Hanwei Xu
Shirong Ma	Honghui Ding
Peiyi Wang	Huazuo Gao
Xiao Bi	Hui Qu
Xiaokang Zhang	Hui Li
Xingkai Yu	Jianzhong Guo
Yu Wu	Jiashi Li
Z.F. Wu	Jingchang Chen
Zhibin Gou	Jingyang Yuan
Zhihong Shao	Jinhao Tu
Zhuoshu Li	Junjie Qiu
Ziyi Gao	Junlong Li
	J.L. Cai
Contributors	Jiaqi Ni
Aixin Liu	Jian Liang
Bing Xue	Jin Chen
Bingxuan Wang	Kai Dong
Bochao Wu	Kai Hu*
Bei Feng	Kaichao You
Chengda Lu	Kaige Gao
Chenggang Zhao	Kang Guan
Chengqi Deng	Kexin Huang
Chong Ruan	Kuai Yu
Damai Dai	Lean Wang
Deli Chen	Lecong Zhang
Dongjie Ji	Liang Zhao
Erhang Li	Litong Wang

Liyue Zhang
Lei Xu
Leyi Xia
Mingchuan Zhang
Minghua Zhang
Minghui Tang
Mingxu Zhou
Meng Li
Miaojun Wang
Mingming Li
Ning Tian
Panpan Huang
Peng Zhang
Qiancheng Wang
Qinyu Chen
Qiusi Du
Ruiqi Ge*
Ruisong Zhang
Ruizhe Pan
Runji Wang
R.J. Chen
R.L. Jin
Ruyi Chen
Shanghao Lu
Shangyan Zhou
Shanhuang Chen
Shengfeng Ye
Shiyu Wang
Shuiping Yu
Shunfeng Zhou
Shuting Pan
S.S. Li
Shuang Zhou
Shaoqing Wu
Shengfeng Ye
Tao Yun
Tian Pei
Tianyu Sun

T. Wang
Wangding Zeng
Wen Liu
Wenfeng Liang
Wenjun Gao
Wenqin Yu*
Wentao Zhang
W.L. Xiao
Wei An
Xiaodong Liu
Xiaohan Wang
Xiaokang Chen
Xiaotao Nie
Xin Cheng
Xin Liu
Xin Xie
Xingchao Liu
Xinyu Yang
Xinyuan Li
Xuecheng Su
Xuheng Lin
X.Q. Li
Xiangyue Jin
Xiaojin Shen
Xiaosha Chen
Xiaowen Sun
Xiaoxiang Wang
Xinnan Song
Xinyi Zhou
Xianzu Wang
Xinxia Shan
Y.K. Li
Y.Q. Wang
Y.X. Wei
Yang Zhang
Yanhong Xu
Yao Li
Yao Zhao

Yaofeng Sun	Yunxian Ma
Yaohui Wang	Ying Tang
Yi Yu	Yukun Zha
Yichao Zhang	Yuting Yan
Yifan Shi	Z.Z. Ren
Yiliang Xiong	Zehui Ren
Ying He	Zhangli Sha
Yishi Piao	Zhe Fu
Yisong Wang	Zhean Xu
Yixuan Tan	Zhenda Xie
Yiyang Ma*	Zhengyan Zhang
Yiyuan Liu	Zhewen Hao
Yongqiang Guo	Zhicheng Ma
Yuan Ou	Zhigang Yan
Yuduan Wang	Zhiyu Wu
Yue Gong	Zihui Gu
Yuheng Zou	Zijia Zhu
Yujia He	Zijun Liu*
Yunfan Xiong	Zilin Li
Yuxiang Luo	Ziwei Xie
Yuxiang You	Ziyang Song
Yuxuan Liu	Zizheng Pan
Yuyang Zhou	Zhen Huang
Y.X. Zhu	Zhipeng Xu
Yanping Huang	Zhongyu Zhang
Yaohui Li	Zhen Zhang
Yi Zheng	
Yuchen Zhu	

Within each role, authors are listed alphabetically by the first name. Names marked with * denote individuals who have departed from our team.