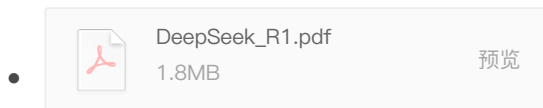


DeepSeek R1 阅读



- DeepSeek R1 zero
 - 基座: DeepSeek v3-base
 - 训练数据
 - 数据量: 未知
 - 训练模版

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

- 训练方法
 - Reward Methods(rule based)
 - **Accuracy rewards:** The accuracy reward model evaluates **whether the response is correct**. For example, in the case of math problems with deterministic results, the model is required to provide the final answer **in a specified format** (e.g., within a box), enabling reliable rule-based verification of correctness. Similarly, for LeetCode problems, a compiler can be used to generate feedback based on predefined test cases.
 - **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.

■ GRPO

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 $\{o_1, o_2, \dots, o_G\}$ from the old policy $\pi_{\theta_{old}}$ and then optimizes the policy model π_{θ} by maximizing the following objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)] \frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta \mathbb{D}_{KL}(\pi_{\theta} || \pi_{ref}) \right), \quad (1)$$

$$\mathbb{D}_{KL}(\pi_{\theta} || \pi_{ref}) = \frac{\pi_{ref}(o_i|q)}{\pi_{\theta}(o_i|q)} - \log \frac{\pi_{ref}(o_i|q)}{\pi_{\theta}(o_i|q)} - 1, \quad (2)$$

where ϵ and β are hyper-parameters, and A_i is the advantage, computed using a group of rewards $\{r_1, r_2, \dots, r_G\}$ corresponding to the outputs within each group:

$$A_i = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}. \quad (3)$$

- DeepSeek R1
 - 基座: DeepSeek v3-base
 - 训练数据
 - 第一阶段数据
 - 数据量: 数千条长CoT数据, 格式规范、语言对齐
 - 来源: DeepSeek R1 Zero蒸馏+清洗
 - 第二阶段数据
 - 数据量: 未知

- 来源：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.

■ 第三阶段数据

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**.

- 数据量：800k

- 推理数据：

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**.

- 非推理数据

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.

- 来源：由第二阶段得到的模型进行拒绝采样+清洗获得

- 清洗方法：拒绝采样+使用DeepSeekv3进行清洗

■ 第四阶段数据：

- 数据量：未知

- 来源：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 simultaneously refining its reasoning capabilities. Specifically, we train the model using a combination of reward signals and diverse prompt distributions.

○ 训练方法

■ 第一阶段（冷启动SFT）

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.

-

- 目的

- 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.

■ 第二阶段（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 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.

- 加入了语言Reward模块 (rule based)

■ 第三阶段 (Rejection Sampling and Supervised Fine-Tuning)

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

■ 第四阶段 (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 simultaneously 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.

□ 采用了两种Reward

- 针对推理数据部分使用rule based

- **Accuracy rewards:** The accuracy reward model evaluates whether the response is correct. For example, in the case of math problems with deterministic results, the model is required to provide the final answer in a specified format (e.g., within a box), enabling reliable rule-based verification of correctness. Similarly, for LeetCode problems, a compiler can be used to generate feedback based on predefined test cases.
- **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.

- 针对非推理数据部分使用DeepSeek V3 (model based)

● 一些值得思考的结论

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 | Comparison of distilled and RL Models on Reasoning-Related Benchmarks.

-
- qwen32B蒸馏比使用RL训练效果好很多

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.