

# FASTCURL: Curriculum Reinforcement Learning with Progressive Context Extension for Efficient Training R1-like Reasoning Models

Mingyang Song, Mao Zheng, Zheng Li, Wenjie Yang, Xuan Luo, Yue Pan, Feng Zhang  
Tencent Hunyuan

nickmysong@tencent.com

## Abstract

In this paper, we propose **FASTCURL**, a simple yet efficient **Curriculum Reinforcement Learning** approach with context window extending strategy to accelerate the reinforcement learning training efficiency for R1-like reasoning models while enhancing their performance in tackling complex reasoning tasks with long chain-of-thought rationales, particularly with a 1.5B parameter language model. **FASTCURL** consists of two main procedures: length-aware training data segmentation and context window extension training. Specifically, the former first splits the original training data into three different levels by the input prompt length, and then the latter leverages segmented training datasets with a progressively increasing context window length to train the reasoning model. Experimental results demonstrate that **FASTCURL-1.5B-Preview** surpasses **DeepScaleR-1.5B-Preview** across all five datasets (including MATH 500, AIME 2024, AMC 2023, Minerva Math, and OlympiadBench) while only utilizing 50% of training steps. Furthermore, all training stages for **FASTCURL-1.5B-Preview** are completed using a single node with 8 GPUs. The code<sup>1</sup>, dataset<sup>2</sup>, and model checkpoints<sup>3</sup> are released.

## 1 Introduction

Large Language Models (LLMs) have emerged as immensely potent AI instruments, showcasing extraordinary proficiency in comprehending natural language and executing downstream tasks (Zhao et al., 2023; Minaee et al., 2024; Tie et al., 2025). Lately, Large Reasoning Models (LRMs), also referred to as slow-thinking models, have made remarkable advancements in strengthening the deliberate and methodical thinking abilities of LLMs, enabling them to tackle complex reasoning tasks with

greater effectiveness (Shao et al., 2024; DeepSeek-AI, 2025; Chen et al., 2025).

Recent groundbreaking developments, such as DeepSeek-R1 (DeepSeek-AI, 2025), have unveiled a scaling phenomenon in the large-scale Reinforcement Learning (RL) training: as more computational resources are allocated to training, there is a persistent and steady increase in both benchmark performance and the length of responses generated by the trained model, with no signs of reaching a saturation point. Inspired by these achievements, training LLMs through large-scale RL has recently emerged as a promising paradigm for addressing complex reasoning tasks. Recently, a wealth of valuable research endeavors has emerged, aiming to explore and replicate reasoning models akin to DeepSeek-R1 (for example, starting from a distilled model or a pre-trained model). Notable examples include DeepScaleR<sup>4</sup>, Open R1<sup>5</sup>, and OpenReasoner-Zero<sup>6</sup>, among others.

In the realm of large-scale reinforcement learning, the huge computational costs stand as one of the biggest challenges (Chen et al., 2025). Recently, DeepScaleR observes that directly replicating DeepSeek-R1’s experiments (with context lengths of at least 32K tokens and roughly 8,000 training steps) requires a minimum of 70,000 A100 GPU hours, even for a relatively small 1.5B parameter language model. To mitigate this issue, DeepScaleR introduces an iterative lengthening strategy for RL, dramatically reducing the computational requirements to 3800 A100 GPU hours. This method conserves resources and outperforms OpenAI’s o1-preview (OpenAI, 2024) using a model with only 1.5 billion parameters. Specifically, DeepScaleR’s reinforcement learning training process involves

<sup>1</sup><https://github.com/nick7nlp/FastCuRL>

<sup>2</sup><https://huggingface.co/datasets/Nickyang/FastCuRL>

<sup>3</sup><https://huggingface.co/Nickyang/FastCuRL-1.5B-Preview>

<sup>4</sup><https://github.com/agentica-project/deepscaler>

<sup>5</sup><https://github.com/huggingface/open-r1>

<sup>6</sup><https://github.com/Open-Reasoner-Zero/Open-Reasoner-Zero>



Figure 1: The first training stage of DeepScaleR.

three distinct stages to enhance the model’s long CoT reasoning capabilities while efficiently utilizing the context window. Initially, analysis of the model’s responses reveals that lengthy outputs often contain repetitive patterns that do not meaningfully contribute to effective long CoT reasoning. To address this, training commenced with an 8K context length to guide the model toward more efficient context use and improve the quality of long CoT rationales. After approximately 1,000 training steps, as the model’s response lengths began to increase again, DeepScaleR extends the training context from 8K to 16K to accommodate this change. Following an additional 500 steps at the 16K context length, when performance improvements start to plateau, the output context length is further increased to 24K, aiming further to enhance the model’s long CoT reasoning abilities.

By analyzing CoT reasoning responses produced by R1-like reasoning models when addressing complex reasoning questions, we have discerned two types of problems that lead to the model generating excessively lengthy responses. The first type pertains to challenging problems requiring long CoT responses to solve. The second involves questions laden with numerous conditions, prompting the model to verify each condition repeatedly during problem-solving. This repetitive verification can result in redundant thinking patterns, ultimately causing the reasoning responses to be unduly long. Both situations may impair the model’s training

efficiency during the 8K context window. Meanwhile, as shown in Figure 1, DeepScaleR exhibits a high clip ratio during the initial training stage. Therefore, segmenting the training data based on specific features (such as response length) could help mitigate truncation due to context limitations, thereby improving the training efficiency.

Inspired by the above observations, we propose **FASTCURL**, a simple yet efficient **Curriculum Reinforcement Learning** approach with context window extending strategy to accelerate the reinforcement learning training efficiency for R1-like reasoning models. Specifically, the former first splits the original dataset into various levels via the input prompt length. Following this, we introduce a curriculum reinforcement learning approach with a progressive increasing context window extension aimed at accelerating the RL training process of R1-like reasoning models, enabling them to master complex reasoning tasks effectively. Experimental results demonstrate that FASTCURL-1.5B-Preview surpasses DeepScaleR-1.5B-Preview across all five popular-used reasoning benchmarks, MATH 500, AIME 2024, AMC 2023, Minerva Math, and OlympiadBench, while reducing 50% training steps compared with DeepScaleR-1.5B-Preview. Furthermore, all training stages for our model FASTCURL-1.5B-Preview are completed using just a single node with 8 GPUs.

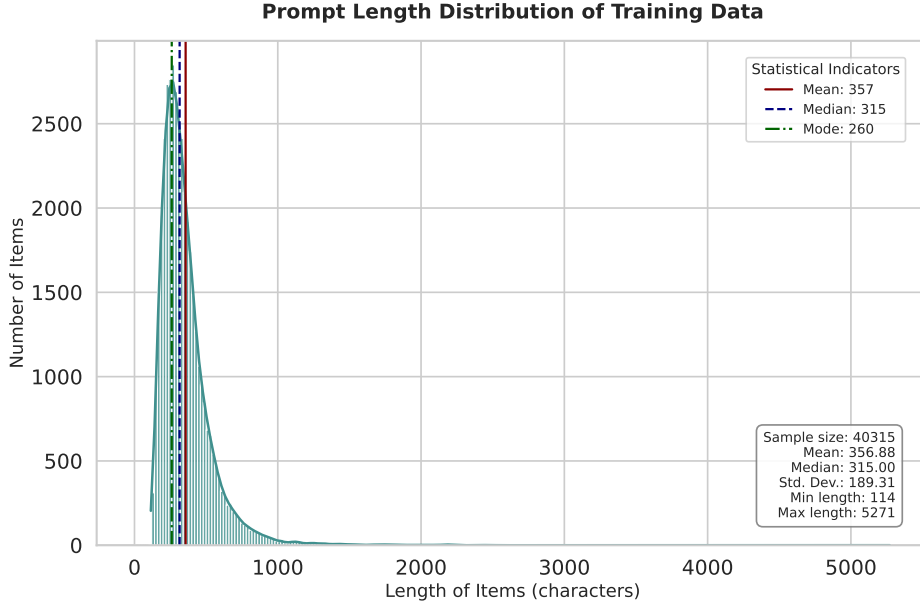


Figure 2: Prompt Length Distribution.

## 2 Methodology

In this section, we describe the details of FastCuRL, which mainly contains two parts: length-sensitive training data segmentation and curriculum reinforcement learning with progressive context window extension.

### 2.1 Length-Sensitive Data Segmentation

As mentioned before, the primary objective of data segmentation is to differentiate between data samples that necessitate lengthy CoT reasoning justifications and those that require only brief CoT rationales. This distinction is crucial for minimizing the adverse effects of truncation when training models within a constrained context window. A seemingly straightforward method would involve running all training data through a reference model and subsequently categorizing the data based on the length of the generated responses. However, this strategy inadvertently imposes constraints from the reference policy onto the training process, thereby influencing its outcome. Additionally, we aim to maintain a degree of uncertainty in the data segmentation process, meaning that the dataset labeled as containing short CoT reasoning examples may include a minor fraction of long CoT reasoning instances and vice versa.

An intuitive approach may employ perplexity, loss functions, or reward models to classify training samples. Nonetheless, the experimental findings presented herein are designed to distinguish

Datasets	Average Input Prompt Length
Short	256
Short+Long	357
Long	521

Table 1: Statistics of Short, Short+Long, Long datasets.

between training data likely to yield lengthy or succinct CoT reasoning responses. This separation aims to facilitate multi-stage reinforcement learning, thereby enhancing training efficiency. Consequently, we propose utilizing more straightforward and resilient criteria for splitting the training data.

As discussed in earlier sections, R1-like reasoning models often engage in profound reflection and validation thoughts when confronted with complex problems featuring numerous conditions (and correspondingly long input prompts). Drawing inspiration from this observation, we postulate that questions with more extensive inputs are generally associated with lengthier outputs. Based on this hypothesis, we statistically analyze the original dataset, as illustrated in Figure 2. Subsequently, we divide the original training dataset (referred to as Short+Long) into two training data subsets: one representing a short CoT reasoning dataset (designated as Short) and the other constituting a long CoT reasoning dataset (labeled as Long) based on the average input length. Finally, the average input length of each dataset as shown in Table 1.

## 2.2 Curriculum Reinforcement Learning with Progressive Context Window Extension

Effectively scaling reinforcement learning for complex reasoning tasks is significantly hampered by the challenge of determining the optimal context window size during training. Complex reasoning tasks inherently require substantial computational resources because they typically generate outputs much longer than standard tasks. This elongation in output length considerably slows down trajectory sampling and policy gradient updates. For instance, doubling the context window length directly results in at least a twofold increase in computational overhead. This situation introduces a fundamental trade-off: longer contexts enhance a model’s cognitive capacity and enable deeper reasoning but come at the cost of substantially lengthier training cycles. Conversely, shorter contexts allow for more efficient training but may limit the model’s ability to tackle complex tasks that require extensive reasoning spans successfully. Therefore, finding an optimal balance between training efficiency and problem-solving accuracy is essential.

As illustrated in Figure 1, over 40% of training samples (may include both lengthy incorrect and correct long CoT reasoning) are truncated at the beginning of training. Therefore, after obtaining the segmented datasets, we train the model with curriculum reinforcement learning by progressively extending the context window. Specifically, we let the model learn step by step: first, the simple things (Short), then the mixed ones (Long+Short), then only the difficult ones (Long), and finally, reviewing everything as a whole once more (Short+Long). Each stage trains at most one iteration of the corresponding dataset. Completing the four stages is equivalent to training for three iterations on the original training data.

## 3 Experimental Settings

In this section, we introduce the used datasets, baselines, training parameters, evaluation setup, and the experimental results.

### 3.1 Datasets

Following DeepScaleR’s methodology, our training dataset consists of 40,315 unique problem-answer pairs collected from the following sources: AIME, AMC, Omni-MATH, and Still (Balunović et al., 2025; Gao et al., 2024; Min et al., 2024). To better evaluate the trained model, we have selected

five benchmarks to assess its performance: MATH 500 (Hendrycks et al., 2021), AIME 2024, AMC 2023, Minerva Math (Lewkowycz et al., 2022), and OlympiadBench (He et al., 2024).

### 3.2 Baselines

In this paper, we conduct comprehensive evaluations against several 1.5B and 7B parameter models as the baselines, which includes DEEPSEEK-R1-DISTILL-QWEN-1.5B, STILL-1.5B<sup>7</sup>, DeepScaleR-1.5B-Preview, RSTAR-MATH-7B (Guan et al., 2025), QWEN-2.5-MATH-7B-Instruct (Yang et al., 2024), QWEN2.5-7B-SimpleRL<sup>8</sup>, and EURUS-2-7B-PRIME (Cui et al., 2025).

### 3.3 Training Parameters

For training parameters, please refer to the GitHub.

### 3.4 Evaluation Setup

Following DeepSeek-AI (2025), we set the maximum generation length for the models to 32,768 tokens and leverage PASS@1 as the evaluation metric. Specifically, we use a sampling temperature of 0.6 and a top-p value of 0.95 to generate  $k$  responses for each question, typically  $k = 16$ . Specifically, PASS@1 is then calculated as:

$$\text{PASS@1} = \frac{1}{k} \sum_{i=1}^k p_i, \quad (1)$$

where  $p_i$  is the correctness of the  $i$ -th response.

### 3.5 Results

Table 2 present the overall PASS@1 performance of QWEN-2.5-MATH-7B-Instruct, DEEPSEEK-R1-DISTILL-QWEN-1.5B, STILL-1.5B, QWEN2.5-7B-SimpleRL, RSTAR-MATH-7B, EURUS-2-7B-PRIME, and DEEPSCALER-1.5B-Preview. Specifically, FASTCURL-1.5B-Preview achieves the best overall performance, which demonstrates the effectiveness of our proposed approach FASTCURL.

Meanwhile, FASTCURL-1.5B-Preview has better generalization on the AMC 2023 and Minerva Math test sets than the baseline DEEPSCALER-1.5B-Preview. Furthermore, as shown in Table 3, compared with the baseline DEEPSCALER-1.5B-Preview, we only use 50% of the training steps during training and only one node with 8 GPUs, saving more than half of the training resources.

<sup>7</sup>[https://github.com/RUCAIBox/Slow\\_Thinking\\_with\\_LLMs](https://github.com/RUCAIBox/Slow_Thinking_with_LLMs)

<sup>8</sup><https://hkust-nlp.notion.site/simplerl-reason>

Model	MATH 500	AIME 2024	AMC 2023	Minerva Math	OlympiadBench	Avg.
QWEN-2.5-MATH-7B-Instruct	79.8	13.3	50.6	34.6	40.7	43.8
RSTAR-MATH-7B	78.4	26.7	47.5	-	47.1	-
EURUS-2-7B-PRIME	79.2	26.7	57.8	38.6	42.1	48.9
QWEN2.5-7B-SimpleRL	82.4	26.7	62.5	<b>39.7</b>	43.3	50.9
DEEPSEEK-R1-DISTILL-QWEN-1.5B	82.8	28.8	62.9	26.5	43.3	48.9
STILL-1.5B	84.4	32.5	66.7	29.0	45.4	51.6
DEEPSALER-1.5B-Preview	87.8	43.1	73.6	30.2	50.0	57.0
<b>FASTCURL-1.5B-Preview</b>	<b>88.0</b>	<b>43.1</b>	<b>74.2</b>	31.6	<b>50.4</b>	<b>57.5</b>

Table 2: PASS@1 accuracy is reported, averaged over 16 samples for each problem. <sup>†</sup> indicates results obtained by re-evaluating using the checkpoints provided by the corresponding work.

Model	Training Steps	Training Stages	Number of GPUs Used in Each Stage
DEEPSALER-1.5B-Preview	~ 1,750	3	8, 16, 32
<b>FASTCURL-1.5B-Preview</b>	~ 860	4	8, 8, 8, 8

Table 3: Training Details. To ensure consistency in counting training steps, we standardized the batch size to 128. This means that two steps with a batch size of 64 are considered equivalent to one step with a batch size of 128.

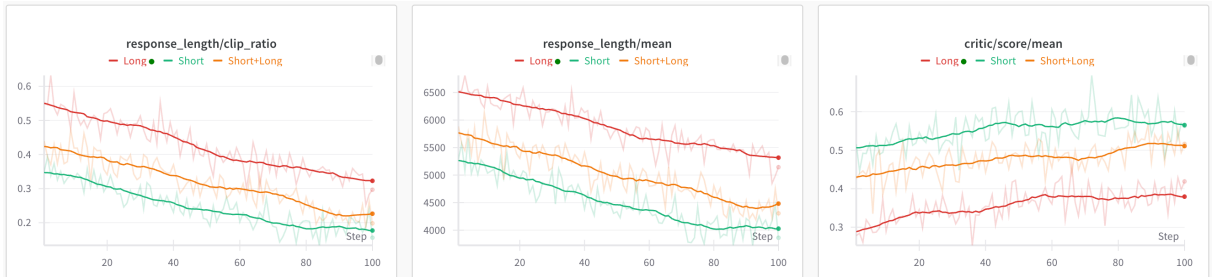


Figure 3: Training records on Short, Short+Long, and Long datasets.

## 4 Discussion

In this section, we analyze the experiments related to segmentation validity verification and the observations of the multiple training processes.

### 4.1 Segmentation Validity Verification

To validate the effectiveness of length-aware data segmentation, we train three models with the same setting on Short, Short+Long, and Long under an 8K context window as seen from Figure 3, whether the experiment results meet expectations in clip ratio, response length, and reward scores. These results demonstrate the effectiveness of the proposed length-aware data segmentation approach.

### 4.2 Training Processes Observation

We have conducted a simple analysis of the four stages in the training process, as follows:

**The First Stage.** Longer responses do not necessarily yield more accurate reasoning results. Therefore, starting training with a long context window might be inefficient, as it can lead to unnecessary token expenditure. We employ the Short dataset and set an 8K context window to enhance training efficiency in the initial stage to optimize the model for generating more concise reasoning rationales.

**The Second Stage.** Figure 4 shows that during the first stage at approximately step 160, the clipping ratio of the model responses drops to around 10%. Subsequently, we set the context window



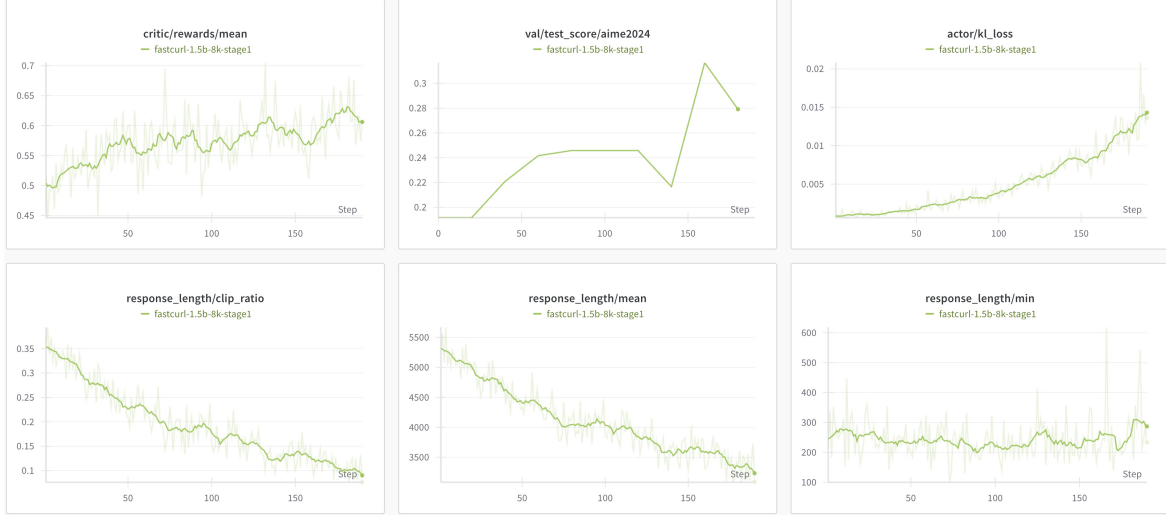


Figure 4: The first training stage log.



Figure 5: The second training stage log.

to 16K and continue training the model using the Short+Long dataset. As shown in Figure 5, the model’s response length and reward have stabilized and are gradually improving. Surprisingly, at step 590, the model’s PASS@1 accuracy on AIME 2024 exceeded 0.4. Thus, one iteration of training for this stage has been completed.

**The Third Stage.** As elaborated earlier, we extend the RL training using the Long dataset in the third stage. Notably, the phenomena observed during this training process closely mirror those seen in the second stage (see Figure 6).

**The Fourth Stage.** In the fourth stage, close to completing one iteration, the model achieves its best performance on AIME 2024 throughout the

entire training process, as shown in Figure 7. With this, the whole training is completed.

Overall, our observations reveal that throughout the training process, the steps selected for stage transitions predominantly occurred near the end of each stage, further emphasizing the effectiveness of the proposed FastCuRL approach. Moreover, we are actively exploring various methods to integrate the Short, Short+Long, and Long datasets. While even faster combination techniques may exist, we are currently in an exploratory phase.

## 5 Conclusions

This paper introduces **FASTCURL**, a straightforward yet highly effective curriculum reinforcement



Figure 6: The third training stage log.

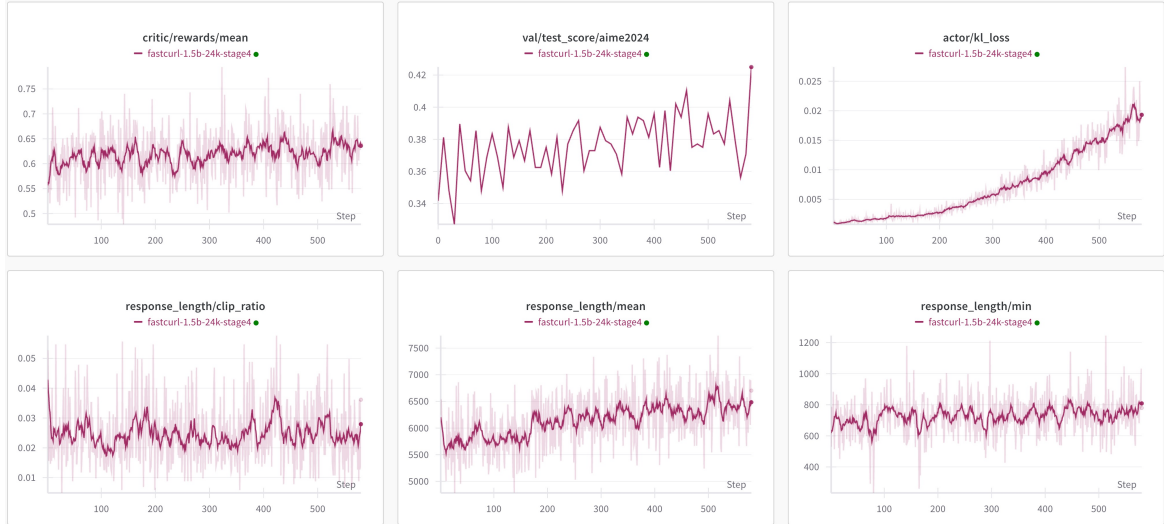


Figure 7: The fourth training stage log.

learning framework incorporating a context extension strategy. This framework is designed to significantly accelerate and enhance the training of R1-like models, especially small language models with approximately 1.5B parameters, in tasks requiring long chains of thought reasoning. Experimental results demonstrate that FASTCURL-1.5B-Preview not only outperforms DeepScaleR-1.5B-Preview across all five benchmarks but also reduces computational resource consumption by more than 50%, with all training phases efficiently executed using just a single node with 8 GPUs.

Training over multiple stages, rather than in a single stage, involves more than changes in parameters like context length; it also fundamentally alters

the reference policy. In a multi-stage training strategy, the KL regularization imposed by the reference policy on the model is gradually relaxed. This reduction in the penalty from the KL divergence between the trained policy and the reference policy allows the trained model to explore a broader range of solutions. Therefore, delving into dynamic control of context window lengths or implementing dynamic KL regularization may be valuable directions for future research.

## References

Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. 2025. [Matharena](#):

- Evaluating llms on uncontaminated math competitions.
- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. [Towards reasoning era: A survey of long chain-of-thought for reasoning large language models](#). *Preprint*, arXiv:2503.09567.
- Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, Jiarui Yuan, Huayu Chen, Kaiyan Zhang, Xingtai Lv, Shuo Wang, Yuan Yao, Xu Han, Hao Peng, Yu Cheng, and 4 others. 2025. [Process reinforcement through implicit rewards](#). *Preprint*, arXiv:2502.01456.
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.
- Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. 2024. Omni-math: A universal olympiad level mathematic benchmark for large language models. *CoRR*, abs/2410.07985.
- Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. 2025. [rstar-math: Small llms can master math reasoning with self-evolved deep thinking](#). *Preprint*, arXiv:2501.04519.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. [Olympiadbench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems](#). In *ACL (1)*, pages 3828–3850. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the MATH dataset](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. [Solving quantitative reasoning problems with language models](#). *Preprint*, arXiv:2206.14858.
- Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, Wayne Xin Zhao, Zheng Liu, Zhongyuan Wang, and Ji-Rong Wen. 2024. [Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems](#). *Preprint*, arXiv:2412.09413.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. [Large language models: A survey](#). *Preprint*, arXiv:2402.06196.
- OpenAI. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *Preprint*, arXiv:2402.03300.
- Guiyao Tie, Zeli Zhao, Dingjie Song, Fuyang Wei, Rong Zhou, Yurou Dai, Wen Yin, Zhejian Yang, Jiangyue Yan, Yao Su, Zhenhan Dai, Yifeng Xie, Yihan Cao, Lichao Sun, Pan Zhou, Lifang He, Hechang Chen, Yu Zhang, Qingsong Wen, and 7 others. 2025. [A survey on post-training of large language models](#). *Preprint*, arXiv:2503.06072.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. [Qwen2.5-math technical report: Toward mathematical expert model via self-improvement](#). *CoRR*, abs/2409.12122.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2023. [A survey of large language models](#). *Preprint*, arXiv:2303.18223.