# Light-R1: Curriculum SFT, DPO and RL for Long COT from Scratch and Beyond

**Liang Wen**[1] **Yunke Cai**[1] **Fenrui Xiao**[1] **Xin He**[1] **Qi An**[1] **Zhenyu Duan**[1]
**Yimin Du**[1] **Junchen Liu**[1] **Lifu Tang**[1] **Xiaowei Lv**[1,2]
**Haosheng Zou**[1] **Yongchao Deng**[1] **Shousheng Jia**[1] **Xiangzheng Zhang**[1]
[1]Qiyuan Tech    [2]Renmin University
zhangxiangzheng@360.cn

## Abstract

This paper presents our work on the Light-R1 series, with models, data, and code all released[1].

We first focus on training long COT models from scratch, specifically starting from models initially lacking long COT capabilities. Using a curriculum training recipe consisting of two-stage SFT and semi-on-policy DPO, we train our model Light-R1-32B from Qwen2.5-32B-Instruct, resulting in superior math performance compared to DeepSeek-R1-Distill-Qwen-32B. Despite being trained exclusively on math data, Light-R1-32B shows strong generalization across other domains. In the subsequent phase of this work, we highlight the significant benefit of the 3k dataset constructed for the second SFT stage on enhancing other models. By fine-tuning DeepSeek-R1-Distilled models using this dataset, we obtain new SOTA models in 7B and 14B, while the 32B model, Light-R1-32B-DS performed comparably to QwQ-32B and DeepSeek-R1.

Furthermore, we extend our work by applying reinforcement learning, specifically GRPO, on long-COT models to further improve reasoning performance. We successfully train our final Light-R1-14B-DS with RL, achieving SOTA performance among 14B models in math. With AIME24 & 25 scores of 74.0 and 60.2 respectively, Light-R1-14B-DS surpasses even many 32B models and DeepSeek-R1-Distill-Llama-70B. Its RL training also exhibits well expected behavior, showing simultaneous increase in response length and reward score.

The Light-R1 series of work validates training long-COT models from scratch, showcases the art in SFT data and releases SOTA models from RL.

Table 1: Light-R1 models. "-DS" = from DeepSeek-R1-Distill, otherwise from Qwen-Instruct.

| Model | AIME24 | AIME25 | GPQA Diamond | Training Recipe |
|---|---|---|---|---|
| Light-R1-32B | 76.6 | 64.6 | 61.8 | SFT stage1&2 + DPO |
| Light-R1-7B-DS | 59.1 | 44.3 | 49.4 | SFT stage2 |
| Light-R1-14B-DS | 74.0 | 60.2 | 61.7 | SFT stage2 + GRPO |
| Light-R1-32B-DS | 78.1 | 65.9 | 68.0 | SFT stage2 |

## 1 Introduction

Since the release of DeepSeek-R1 [DeepSeek-AI, 2025], long chain-of-thought (long COT) reasoning has gained widespread popularity in both foundational AI models and various industrial AI applica-

---

[1]https://github.com/Qihoo360/Light-R1

tions. However, deploying full-capacity R1-level models (typically 70B+ parameters, DeepSeek-R1 with 671B parameters) incurs prohibitive computational costs [DeepSeek-AI, 2025, Qwen, 2025]. The resource barrier of training and deploying the giant models makes them impractical for edge devices and real-time applications. This limitation has sparked growing interest in developing compact yet capable models under a few 10B parameters that can perform extended long COT - a critical requirement for mathematical problem solving, algorithmic planning, and scientific analysis [Wei et al., 2022].

To this end, we conduct this work on the Light-R1 series. At the origin of everything, we established stable and trustworthy evaluation protocols that could precisely reproduce the evaluation results as reported in DeepSeek-AI [2025].

Based on trustworthy evaluation, our research addresses three fundamental challenges in this direction through systematic architectural and algorithmic innovations.

The first challenge lies in curating an efficient and effective dataset for Post-Training. A meticulously designed data selection strategy constitutes the cornerstone of all potential improvements in this stage. To address this challenge, we collected a diverse set of open-source reasoning datasets spanning mathematical reasoning, logical deduction and algorithmic problem-solving. The raw dataset underwent rigorous preprocessing to eliminate duplicates and standardize formatting. Subsequently, we implemented a sophisticated two-stage difficulty filtering methodology to identify the most valuable training examples. Specifically, we employed a sequential evaluation approach utilizing DeepScaleR-1.5B-Preview and DeepSeek-R1-Distill-Qwen-32B models to quantify question difficulty based on pass rate metrics.

The second challenge then emerges as how to optimize the utilization of this dataset. While conventional approaches typically employ a single SFT stage, our preliminary experiments revealed significant limitations with this method for long reasoning problems. Specifically, after the initial SFT stage with our 32B model, we observed that nearly 20% of the training data still exhibited pass rates below 50% across 10 runs, indicating that a single training phase is insufficient to fully assimilate the knowledge contained in datasets with heterogeneous difficulty levels. To address this limitation, we implemented a curriculum learning strategy to maximize the dataset's value. Our experiments demonstrated that the optimal post-training curriculum varies depending on the base model's characteristics—models without inherent long reasoning capabilities typically require more training stages than those already possessing such abilities. For our Light-R1-32B model, which was trained from scratch, optimal performance necessitated a curriculum of two consecutive stages of SFT with increasing difficulty levels, followed by a DPO stage.

The third challenge arises from implementing the final component of Post-Training—Reinforcement Learning (RL)—to further enhance model performance. We are excited to report our successful reinforcement learning (RL) training of Light-R1-14B-DS. Recent RL works have successfully trained RL on base or short-COT models (usually with -zero in their names) [Zeng et al., 2025, Hu et al., 2025], or on small models (with response length interestingly decreasing significantly then increasing) [Zeng et al., 2025, Luo et al., 2025], or on QwQ-32B with presumably prohibitively heavy compute. Among them, none except QwQ-32B could be regarded as full-scale RL reproduction of the non-zero RL of DeepSeek-R1, while QwQ-32B performed RL on 32B models with "scaling RL" [Qwen, 2025] - presumably prohibitively heavy compute. Our long-COT RL Post-Training is the first to show simultaneous increase in response length and reward score on long-COT 14B models with no length dropping in the beginning. This breakthrough demonstrates that careful curriculum design can overcome the known scalability limitations of RL in small models Gao et al. [2023].

The key contributions of this work include:

- A detailed, fully open-source Post-Training approach to train long-COT models from scratch. Curriculum SFT + DPO is validated on Qwen2.5-32B-Instruct and could be easily migrated to Qwen2.5-7B and 14B models. The total three-stage curriculum (SFT→SFT→DPO) incrementally builds reasoning capacity through difficulty-progressive data exposure, requiring only $1000 training cost (6 hours on 12×H800 GPUs).

- A well established SFT stage 2 dataset of 3k mostly math questions that could significantly improve not only SFT stage 1 but also all DeepSeek-R1-Distill models, resulting in our SOTA 7B model Light-R1-7B-DS.

Table 2: Reproduction of DeepSeek-AI [2025] and Qwen [2025] evaluation results

| Model | AIME24 pass@1, paper | AIME24 pass@1, ours |
|---|---|---|
| DeepSeek-R1-distill-Qwen-32B | 72.6 | 72.3 |
| DeepSeek-R1-distill-Qwen-14B | 69.7 | 69.3 |
| DeepSeek-R1-distill-Qwen-7B | 55.5 | 54.0 |
| QWQ-32B | 79.5 | 78.5 |

- First demonstration of RL effectiveness on 14B models for mathematical reasoning, achieving around 2% absolute improvement compared with before-RL, resulting in our SOTA 14B model Light-R1-14B-DS.

Our resulting Light-R1-32B model achieves 76.6% on AIME24 and 64.6% on AIME25, surpassing DeepSeek-R1-Distill-Qwen-32B by 4.0% and 9.7%. Light-R1-7B-DS and Light-R1-14B-DS are SOTA math models of equal sizes. Light-R1-14B-DS shows consistent and expected improvement during RL training.

The Light-R1 model series establishes new possibilities for deploying advanced reasoning capabilities in resource-constrained environments such as edge computing without sacrificing analytical depth.

## 2 The Origin of Everything: Stable and Trustworthy Evaluation of Long-COT Models

Following DeepSeek-AI [2025], long-COT models are commonly deployed with sampling temperature 0.6. While long-COT models generally perform better with sampling than with greedy decoding, it brings more burden for model evaluation as multiple samples for each question may be required, contrary to previous viable approaches of greedy decoding for evaluation [Song et al., 2024].

DeepSeek-AI [2025] generates 64 responses per query to estimate pass@1. We have verified this choice, witnessing large deviation of over 3 points using 16 responses or fewer across different runs of the same model. Such randomness is unacceptable to compare model performances.

For stable and trustworthy evaluation, we adapted Luo et al. [2025]'s evaluation code for all our evaluation runs. Our evaluation code and logs are all released.

We can reproduce all DeepSeek-R1-Distil models' and QwQ's scores as reported in DeepSeek-AI [2025], Qwen [2025] as shown in Tab. 2 with 64 samples per query, with deviation around 1 point. We could therefore train and evaluate our models in a trusted way using the same evaluation code and protocol.

## 3 Light-R1-32B: Long COT from Scratch with Curriculum SFT & DPO

While much work [Ye et al., 2025, Muennighoff et al., 2025, Team, 2025] has been open-sourced trying to reproduce DeepSeek-R1 on models of 72B or less, none achieves similar performance on the hard math competition AIME24 & 25 as DeepSeek-R1-Distill-Qwen-32B's score 72.6 & 54.9.

Light-R1-32B starts from Qwen2.5-32B-Instruct without long COT (*from scratch* in terms of long-COT) and trains on decontaminated math data. It distills DeepSeek-R1 with curriculum SFT & DPO to surpass DeepSeek-R1-Distill-Qwen-32B on AIME24 & 25, and improved further with model merging, achieving 76.6 & 64.6.

More importantly, besides the state-of-the-art from-scratch model Light-R1-32B, we have also released all training datasets of our curriculum SFT & DPO and training code based on 360-LLaMA-Factory [Zou et al., 2024]. Estimated training time on 12 x H800 machines takes no more than 6 hours — around $1000.

We believe Light-R1 represents a practical way of training strong long COT models from scratch (from models without long COT). While we have made RL work on 14B models (Sec. 4), curriculum SFT & DPO facilitates more control along the pipeline and is more cost-friendly.

We introduce our data processing and Post-Training pipeline in this section as illustrated by Fig. 1.
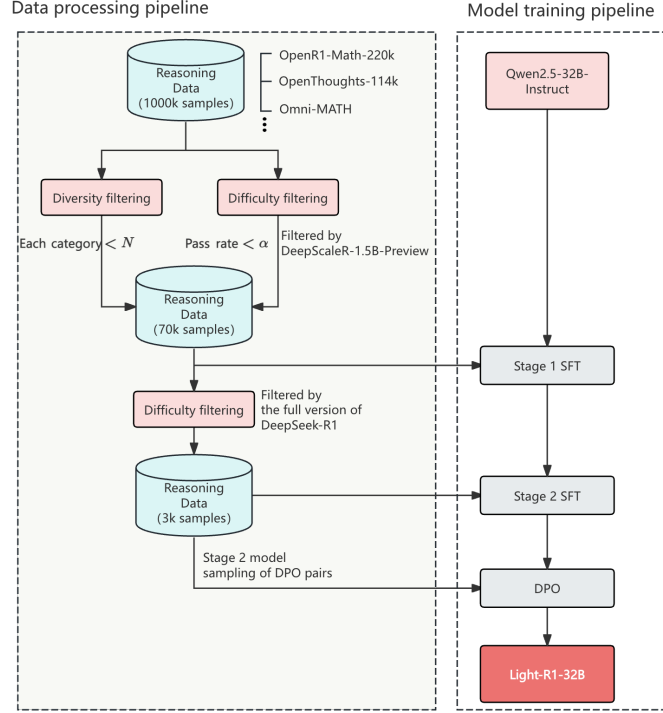


Figure 1: Overview of training pipeline of Light-R1-32B.

## 3.1 Data Preparation

The whole data preparation process spans data collection, data decontamination and data generation, detailed as follows.

### 3.1.1 Data Collection

We began by collecting various sources of math questions with groundtruth answers. Iterating over all possible sources by the time, we collected around 1000k math questions as the seed set. Major data sources include:

- `https://huggingface.co/datasets/GAIR/LIMO`
- `https://huggingface.co/datasets/nvidia/OpenMathInstruct-2`
- `https://huggingface.co/datasets/open-r1/OpenR1-Math-220k`
- `https://huggingface.co/datasets/open-thoughts/OpenThoughts-114k`
- `https://huggingface.co/datasets/simplescaling/s1K-1.1`
- `https://huggingface.co/datasets/KbsdJames/Omni-MATH`
- `https://hf-mirror.com/datasets/baber/hendrycks_math`

All data are aggregated together to form around 1000k math questions as the seed set. Within this 1000k data, we kept only math questions with groundtruth answers. Questions without groundtruth answers could be used as synthetic data by letting multiple strong LLMs vote for groundtruths but we left it for future work.

The data is then filtered for diversity, where we tagged each question with an in-house tagging system and downsample categories with excessive data.

### 3.1.2 Data Decontamination

Table 3: Number of matched prompts in open-source datasets against benchmarks.

| Dataset | AIME24+25 | MATH-500 | GPQA Diamond |
|---|---|---|---|
| OpenThoughts-114k | 0 | 100 | 0 |
| Open-R1-Math-220k | 0 | 10 | 0 |
| DeepScaleR-Preview-Dataset | 0 | 196 | 0 |
| LIMO | 0 | 0 | 0 |
| Bespoke-Stratos-17k | 0 | 125 | 0 |
| Open-Reasoner-Zero | 0 | 325 | 0 |
| simplescaling/data_ablation_full59K | 0 | 244 | 1 |
| simplescaling/s1K-1.1 | 0 | 3 | 1 |
| ours | 0 | 0 | 0 |

We carefully evaluated data contamination of several open-sourced datasets. While certain contamination may be inevitable during pre-training, it is unacceptable for post-training to compare on benchmarks. As can be seen in Tab. 3, we found that MATH-500 is somewhat compromised with tens of questions that are identical or only numbers changed. AIME 24 and 25 stay intact but we have to pay special attention when we incorporate AIME data up to 2023.

Light-R1 did thorough decontamination with exact matching (excluding digits, to filter out questions with only digits changed) and N-gram (N=32) matching against AIME24, AIME25, MATH-500 and GPQA.

### 3.1.3 Data Generation

With a diverse and clean dataset, we generate long-COT responses for SFT training. However, not all data are suitable and necessary to train on, and distilling DeepSeek-R1 could be expensive whether querying API or deploying locally. We therefore performed difficulty first on the dataset, to keep only questions that are not too easy.

We use Luo et al. [2025]'s DeepScaleR-1.5B-Preview model to sample responses for each question, as the model is small but strong enough. Only questions with pass rate $< \alpha$ are kept to query DeepSeek-R1. This result in around 70k (76k to be precise) data. After querying DeepSeek-R1, only questions with correct long-COT answers are kept. If two or more sampled answers are correct, we randomly choose one of the long-COT answers for SFT.

This way, we constructed an over 70k SFT dataset, whose prompts are filtered by diversity and difficulty and long-COT responses are generated by DeepSeek-R1 and verified against the groundtruth.

However, directly training on this dataset cannot produce satisfactory results directly no matter how many epochs SFT is trained. After inspecting the trained model' performance on different questions, we found that the model may need further training on more difficult questions. Therefore, instead of difficulty filtering with DeepScaleR-1.5B-Preview in the last stage, we performed another stage of difficulty filtering with DeepSeek-R1-Distill-Qwen-32B. Only questions with pass rate $< \alpha$ and questions which DeepSeek-R1's sampled responses cannot get all right or all wrong are kept in this stage, resulting in a SFT stage 2 dataset of size 3k. Interestingly, this dataset exhibits such high quality that training solely on it yields performance improvements across all DeepSeek-R1-Distill models, as will be discussed in Section 3.4.

### 3.2 Curriculum Post-Training

Our approach consists of three stages:

1. **SFT Stage 1**: Training on 76k filtered mathematical problems
2. **SFT Stage 2**: Fine-tuning on 3k high-difficulty problems

Table 4: Stage-wise performance improvement

| Stage | AIME24 | AIME25 | GPQA Diamond |
|---|---|---|---|
| Qwen2.5-32B-Instruct (base model) | 16.6 | 13.6 | 48.8 |
| Light-R1-32B-SFT-stage1 | 69.0 | 57.4 | 64.3 |
| Light-R1-32B-SFT-stage2 | 73.0 | 64.3 | 60.6 |
| Light-R1-32B-DPO | 75.8 | 63.4 | 61.8 |
| Light-R1-32B (merged model) | **76.6** | **64.6** | 61.8 |

Table 5: Improvement with 3k data.

| Model | AIME24 | AIME25 | GPQA Diamond |
|---|---|---|---|
| DeepSeek-R1-Distill-Qwen-7B | 55.5 | 39.2 | 49.1 |
| Light-R1-7B-DS | **59.1** | **44.3** | 49.4 |
| DeepSeek-R1-Distill-Qwen-14B | 69.7 | 50.2 | 59.1 |
| Light-R1-14B-DS-SFT | **72.3** | **58.9** | N/A |
| DeepSeek-R1-Distill-Qwen-32B | 72.6 | 54.9 | 62.1 |
| Light-R1-32B-DS | **78.1** | **65.9** | **68.0** |

3. **DPO Optimization**: Preference-based optimization using verified response pairs

Both SFT stages are trained with the two-stage data as discussed in Sec. 3.1.3. Detailed hyperparameters could be found in our GitHub repo.

For DPO we adopted a semi-on-policy method with NCA loss [Chen et al., 2024]. Rejected responses are sampled from the SFT-stage-2 model with verified wrong answer. Specifically, responses too long, too short or think right but answer wrong are selected as rejected in DPO pairs. Chosen responses are verified correct answers from DeepSeek-R1 While we have been using fully on-policy DPO for a long time, we found that for hard math problems the chosen responses are better from much stronger models.

## 3.3 Results

We have seen steady improvement in each stage of our curriculum SFT & DPO Post-Training, as shown in Tab. 4. Following DPO stage, we merged models of SFT-stage2, DPO and another DPO version with AIME24 score 74.7. The two DPO versions differ in that one of the data has special tokens skipped in rejected responses. Interestingly, the resulting version also exhibits improvement. On the GPQA evaluation of scientific questions we didn't train on at all, math-specialized training has led to some degree of forgetting. However, Light-R1-32B still demonstrates strong generalization ability.

## 3.4 High-Quality Data is All You Need

Considering DeepSeek-R1-Distill-Qwen models as a stronger version of our SFT stage 1, we performed SFT stage 2 with the 3k stage 2 data on top of DeepSeek-R1-Distill-Qwen models.

Surprisingly as Tab. 5, we could achieve universal improvement on DeepSeek-R1-Distill-Qwen models with this 3k data alone, demonstrating the high quality of the stage 2 data. It may also be because this 3k data is to some extent orthogonal to DeepSeek-R1-Distill-Qwen models' 800k SFT data, hence such easy improvement.

Light-R1-7B-DS and Light-R1-32B-DS are released as the directly SFT stage 2 version, while Light-R1-14B-DS has undergone additional RL training for further improvement.

GPQA is unexpectedly high for Light-R1-32B-DS, but in most of our experiments science and coding other than math should be further improved by specific training, which we haven't done in this technical report.
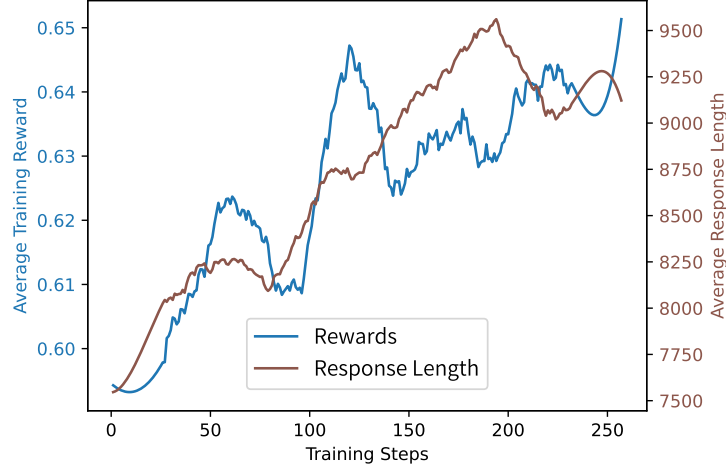
Figure 2: RL Learning curves of response length and train-reward, smoothed with Savitzky-Golay filter.

Table 6: RL performance improvement of Light-R1-14B-DS.

| Model | AIME24 | AIME25 | GPQA Diamond |
| --- | --- | --- | --- |
| DeepSeek-R1-Distill-Qwen-14B | 69.7 | 50.2 | 59.1 |
| Light-R1-14B-DS-SFT | 72.3 | 58.9 | N/A |
| Light-R1-14B-DS-SFT-GPRO epoch1 | 72.3 | 57.8 | N/A |
| Light-R1-14B-DS-SFT-GPRO epoch2 | 73.4 | 60.5 | N/A |
| Light-R1-14B-DS(-SFT-GPRO epoch3) | **74.0** | **60.2** | **61.7** |

# 4 Light-R1-14B-DS: Successful RL on Already Long-COT Finetuned Models

We conduct our reinforcement learning experiments on DeepSeek-R1-Distill-Qwen-14B. To the best of our knowledge, this is the first publicly documented work demonstrating significant improvement in performance through RL on already long-COT 14B models.

Previous studies by DeepSeek-AI [2025], Yuan et al. [2025], and Zhang et al. [2025] have shown that smaller models (with 32 billion parameters or fewer) can reach high performance levels through distillation from larger reasoning models. However, further improvement via RL (Reinforcement Learning) on already long-COT finetuned models is not yet widely reached by the community and is not as easily reachable as *zero* RL (Sec. 1). While Luo et al. [2025] successfully demonstrated promising RL training on a smaller model DeepSeek-R1-Distill-Qwen-1.5B, we encountered challenges in replicating similar results with the larger DeepSeek-R1-Distill-Qwen-14B model using the same recipe.

After weeks of investigation, we arrived at our final RL solution consisting of a two-step process, drawing inspiration from our effective curriculum SFT attempt and Cui et al.. The process is as follows:

1. **Offline Data Selection**: Use Light-R1-7B-DS to sample results of RL training prompts. Keep only the prompts whose pass rates are not 0 or 1 and within a certain range.

2. **Online Reinforcement Learning**: Apply GRPO on the filtered dataset.

We choose GRPO ([Shao et al.]) as the optimization algorithm and implement it based on verl ([Sheng et al., 2024]). We also employ two techniques to stabilize the RL training process: modified version of length reward [Yeo et al.] with weaker preference for short correct answers and importance sampling weight clipping [MiniMax et al.].

We use a rule-based reward and the de-duplicated version of the Big-Math dataset (Albalak et al. [2025]). The experiments are conducted on a cluster of 16 * 8 A100 GPUs. The offline data selection

process takes 4 hours, while the online reinforcement learning takes 26 hours to complete 140 steps and 42 hours to complete 220 steps.

As can be seen from Fig. 2, our RL training demonstrates expected behavior: simultaneous increase in response length and reward score. No interesting length dropping in the beginning. We evaluated RL epochs 1 and 2 after we finished training 3 epochs. As shown in Tab. 6, although first two epochs seem to bring not much improvement, the healthy RL training curves offer us confidence to continue training. Light-R1-14B-DS is finally RL trained for around 3 epochs, or 220 steps.

## 5  Conclusion

Our Light-R1 series systematically addresses the challenge of training long-chain-of-thought (COT) mathematical reasoning models under resource constraints, offering three significant contributions. First, we establish a reproducible and cost-effective curriculum using SFT and DPO that successfully develops long-COT capabilities from scratch. Second, our carefully curated 3K stage-2 SFT dataset demonstrates remarkable transferability across various model sizes and architectures, significantly enhancing DeepSeek-R1-Distill models and establishing new performance benchmarks for models with 7B, 14B, and 32B parameters. Third, we present the first successful Reinforcement Learning training implementation on a 14B parameter long-COT model, Light-R1-14B-DS, which achieves superior performance while maintaining stable response length growth throughout the training process.

These advancements not only democratize access to R1-level reasoning capabilities but also provide valuable insights into curriculum design, data efficiency, and RL scalability for long reasoning models. Our open-source models, datasets, and code (available at GitHub: Qihoo360/Light-R1) aim to accelerate research in developing compact yet powerful reasoning systems, particularly for resource-constrained applications. Future work will explore the integration of enhanced generalization capabilities for long reasoning models and further optimization of RL training efficiency.

## References

DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL `https://arxiv.org/abs/2501.12948`.

Qwen. Qwq-32b: Embracing the power of reinforcement learning, 2025. URL `https://qwenlm.github.io/blog/qwq-32b/`.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Weihao Zeng, Yuzhen Huang, Wei Liu, Keqing He, Qian Liu, Zejun Ma, and Junxian He. 7b model and 8k examples: Emerging reasoning with reinforcement learning is both effective and efficient. `https://hkust-nlp.notion.site/simplerl-reason`, 2025. Notion Blog.

Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, and Heung-Yeung Shum Xiangyu Zhang. Open-reasoner-zero: An open source approach to scaling reinforcement learning on the base model. `https://github.com/Open-Reasoner-Zero/Open-Reasoner-Zero`, 2025.

Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. Notion Blog.

Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pages 10835–10866. PMLR, 2023.

Yifan Song, Guoyin Wang, Sujian Li, and Bill Yuchen Lin. The good, the bad, and the greedy: Evaluation of llms should not ignore non-determinism. *arXiv preprint arXiv:2407.10457*, 2024.

Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning, 2025. URL `https://arxiv.org/abs/2502.03387`.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

OpenThoughts Team. Open Thoughts. https://open-thoughts.ai, January 2025.

Haosheng Zou, Xiaowei Lv, Shousheng Jia, and Xiangzheng Zhang. 360-llama-factory, 2024. URL `https://github.com/Qihoo360/360-LLaMA-Factory`.

Huayu Chen, Guande He, Hang Su, and Jun Zhu. Noise contrastive alignment of language models with explicit rewards. *arXiv preprint arXiv:2402.05369*, 2024.

Yufeng Yuan, Yu Yue, Ruofei Zhu, Tiantian Fan, and Lin Yan. What's Behind PPO's Collapse in Long-CoT? Value Optimization Holds the Secret, March 2025.

Hanning Zhang, Jiarui Yao, Chenlu Ye, Wei Xiong, and Tong Zhang. Online-dpo-r1: Unlocking effective reasoning without the ppo overhead, 2025. Notion Blog.

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, Zhiyuan Liu, Maosong Sun, Bowen Zhou, and Ning Ding. Process Reinforcement through Implicit Rewards. URL `http://arxiv.org/abs/2502.01456`.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. URL `http://arxiv.org/abs/2402.03300`.

Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.

Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying Long Chain-of-Thought Reasoning in LLMs. URL `http://arxiv.org/abs/2502.03373`.

MiniMax, Aonian Li, Bangwei Gong, Bo Yang, Boji Shan, Chang Liu, Cheng Zhu, Chunhao Zhang, Congchao Guo, Da Chen, Dong Li, Enwei Jiao, Gengxin Li, Guojun Zhang, Haohai Sun, Houze Dong, Jiadai Zhu, Jiaqi Zhuang, Jiayuan Song, Jin Zhu, Jingtao Han, Jingyang Li, Junbin Xie, Junhao Xu, Junjie Yan, Kaishun Zhang, Kecheng Xiao, Kexi Kang, Le Han, Leyang Wang, Lianfei Yu, Liheng Feng, Lin Zheng, Linbo Chai, Long Xing, Meizhi Ju, Mingyuan Chi, Mozhi Zhang, Peikai Huang, Pengcheng Niu, Pengfei Li, Pengyu Zhao, Qi Yang, Qidi Xu, Qiexiang Wang, Qin Wang, Qiuhui Li, Ruitao Leng, Shengmin Shi, Shuqi Yu, Sichen Li, Songquan Zhu, Tao Huang, Tianrun Liang, Weigao Sun, Weixuan Sun, Weiyu Cheng, Wenkai Li, Xiangjun Song, Xiao Su, Xiaodong Han, Xinjie Zhang, Xinzhu Hou, Xu Min, Xun Zou, Xuyang Shen, Yan Gong, Yingjie Zhu, Yipeng Zhou, Yiran Zhong, Yongyi Hu, Yuanxiang Fan, Yue Yu, Yufeng Yang, Yuhao Li, Yunan Huang, Yunji Li, Yunpeng Huang, Yunzhi Xu, Yuxin Mao, Zehan Li, Zekang Li, Zewei Tao, Zewen Ying, Zhaoyang Cong, Zhen Qin, Zhenhua Fan, Zhihang Yu, Zhuo Jiang, and Zijia Wu. MiniMax-01: Scaling Foundation Models with Lightning Attention. URL `http://arxiv.org/abs/2501.08313`.

Alon Albalak, Duy Phung, Nathan Lile, Rafael Rafailov, Kanishk Gandhi, Louis Castricato, Anikait Singh, Chase Blagden, Violet Xiang, Dakota Mahan, and Nick Haber. Big-math: A large-scale, high-quality math dataset for reinforcement learning in language models, 2025. URL `https://arxiv.org/abs/2502.17387`.