# Learning Generative Selection for Best-of-N

**Shubham Toshniwal, Aleksander Ficek*, Siddhartha Jain*, Wei Du, Vahid Noroozi, Sadegh Mahdavi, Somshubra Majumdar, Igor Gitman**

**Abstract:** Scaling test-time compute via parallel sampling can substantially improve LLM reasoning, but is often limited by Best-of-$N$ selection quality. Generative selection methods, such as GenSelect [20], address this bottleneck, yet strong selection performance remains largely limited to large models. We show that small reasoning models can acquire strong GenSelect capabilities through targeted reinforcement learning. To this end, we synthesize selection tasks from large-scale math and code instruction datasets by filtering to instances with both correct and incorrect candidate solutions, and train 1.7B-parameter models with DAPO [26] to reward correct selections. Across math (AIME24, AIME25, HMMT25) and code (LiveCodeBench) reasoning benchmarks, our models consistently outperform prompting and majority-voting baselines, often approaching or exceeding much larger models. Moreover, these gains generalize to selecting outputs from stronger models despite training only on outputs from weaker models. Overall, our results establish reinforcement learning as a scalable way to unlock strong generative selection in small models, enabling efficient test-time scaling.

## 1. Introduction

Test-time scaling has emerged as a powerful approach for improving the reasoning performance of large language models (LLMs) [15, 5]. A simple and widely used strategy is *parallel sampling*: generating multiple independent responses and selecting the best [2, 19]. In practice, the effectiveness of parallel sampling is often limited by the quality of Best-of-$N$ selection (i.e., choosing the best answer among $N$ candidates), except for domains with automatic verifiers.

Existing Best-of-$N$ methods typically rely on shallow aggregation mechanisms, such as majority voting or self-consistency [23], or on external reward models that score each candidate independently [4, 8]. Recently, a new paradigm has emerged that uses reasoning models to perform aggregation over the $N$ responses, posing this aggregation as a reasoning problem [31, 16]. Among these, `GenSelect` [21] prompts a reasoning model to explicitly select the best candidate (given $N$ candidate solutions, output the index of the selected candidate) and has demonstrated strong gains over prior methods.

In this work, we ask whether strong GenSelect behavior can be learned by *small* reasoning models. We propose a targeted reinforcement learning framework for generative selection on math and coding tasks, where models are rewarded for selecting a correct solution among imperfect candidates. We construct candidate sets from problems on which the candidate generator achieves $\leq 50\%$ pass rate, ensuring nontrivial selection difficulty, and determine correctness using automatic verifiers. Using on-policy reinforcement learning, we train 1.7B-parameter reasoning models that consistently outperform prompting and majority-voting baselines and approach the selection performance of much larger 8B-scale models ($\sim 5\times$ larger). Notably, these gains generalize to selecting higher-quality outputs produced by stronger generation models, suggesting that learned selection transfers beyond the training-time candidate distribution.

Our contributions are:

- We introduce a reinforcement learning framework that trains small reasoning models to perform generative Best-of-$N$ selection.
- We construct math and coding selection tasks with automatically verified correctness and controllable selection difficulty.
- We show that 1.7B selectors substantially improve Best-of-$N$ performance and transfer to selecting candidates from stronger generators.

## 2. Data Construction

We construct GenSelect-style selection tasks for math and code using `Qwen3-1.7B` to generate candidate solutions and automatic verifiers (symbolic math checking or unit tests) to label correctness. We filter prompts to ensure they contain at least one correct solution while remaining nontrivial.

### 2.1. Math

We construct math selection data from `OpenMathReasoning` [13] by selecting 37K problems on which the base model (`Qwen3-1.7B`) achieves a pass rate below 50% under a symbolic verifier based on Math-Verify [9].

---

*Equal contribution.
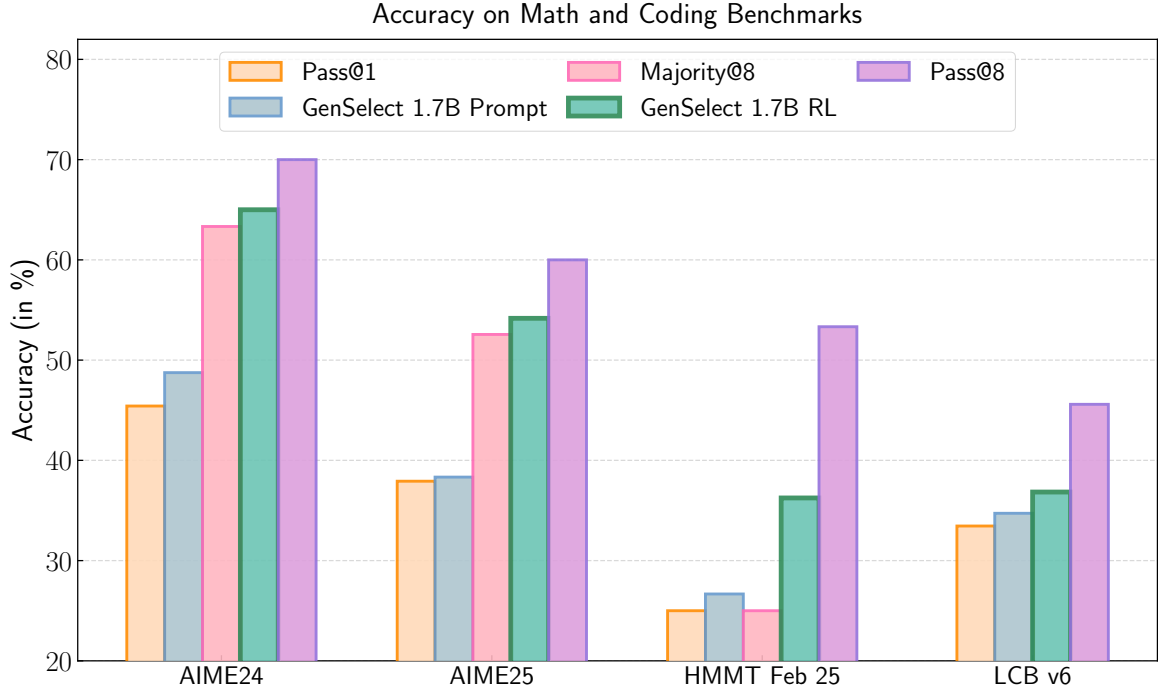
Accuracy on Math and Coding Benchmarks



Figure 1: Accuracy of `Qwen3-1.7B` generations when selected using different strategies on math and coding benchmarks; RL-trained GenSelect consistently outperforms prompting and majority voting.

---

**GenSelect Prompt**

```
You will be given a problem followed by
an enumerated list of {num_solutions}
candidate solutions.  Your task is to
systematically analyze these solutions to
identify the best approach.

Problem:  {problem}

Solutions:  {solutions}

End your evaluation with exactly:

Judgment [IDX]

where IDX is the index 0-{max_idx} of the
best solution.
```

Figure 2: The GenSelect prompt used for selection in prompt baselines and RL training.

For each problem, we sample 2–16 candidate solutions from `Qwen3-1.7B` and assemble them into GenSelect-style prompts (Figure 2). We require at least one correct solution per prompt and cap the fraction of correct candidates at 50% to keep the selection task nontrivial. We construct up to four prompts per problem by resampling candidate solution sets using the above recipe; for training efficiency, we discard prompts exceeding 16K tokens.

## 2.2. Code

For code, we use problems and test cases from `OpenCodeReasoning` (OCR) [1]. We generate candidate solutions with `Qwen3-1.7B` and execute them against the provided unit tests to obtain binary correctness labels.

We then construct GenSelect prompts analogously to the math setting (Figure 2), including up to 16 candidate solutions from `Qwen3-1.7B` and enforcing at least one correct solution with a 50% correctness cap. We do this for both domains to bias the dataset toward harder problems and balance candidate sets to contain both correct and incorrect solutions, making the selection task nontrivial. We discard prompts exceeding 16K tokens and responses over 12K tokens.

## 3. Experimental Setup

### 3.1. RL Training

We train `Qwen3-1.7B` models for both math and code using the VeRL framework [17] on NVIDIA H100 GPUs. We use the DAPO algorithm [26] with on-policy selection generation. We use a batch size of 128, 16 rollouts for math and 8 rollouts for code, a learning rate of $1 \times 10^{-6}$, and the AdamW optimizer. For sampling, we use temperature 1.5, top-$p$ 1.0, and a max output length of 16,384 tokens for both models.

|  | AIME24 | AIME25 | HMMT25 | LCB v6 |
|---|---|---|---|---|
| **Generator: Qwen3-1.7B** | | | | |
| Pass@1 | $45.42 \pm 5.33$ | $37.92 \pm 6.16$ | $25.00 \pm 7.13$ | $33.45 \pm 1.08$ |
| Pass@8 | 70.00 | 60.00 | 53.33 | 45.59 |
| Majority@8 | 63.33 | 52.56 | 25.00 | N/A |
| GenSelect 1.7B Prompting | $48.75 \pm 6.89$ | $38.33 \pm 5.04$ | $26.67 \pm 5.63$ | $34.72 \pm 0.71$ |
| GenSelect 4B Prompting | $61.25 \pm 4.34$ | $53.75 \pm 4.15$ | $34.17 \pm 2.95$ | $39.70 \pm 0.81$ |
| GenSelect 8B Prompting | $\underline{65.83} \pm 2.36$ | $\underline{55.42} \pm 3.05$ | $\underline{38.75} \pm 3.05$ | $\underline{41.24} \pm 0.67$ |
| GenSelect 1.7B RL Math | $\mathbf{65.00} \pm 3.98$ | $\mathbf{54.17} \pm 2.36$ | $\mathbf{36.25} \pm 4.15$ | $36.45 \pm 0.86$ |
| GenSelect 1.7B RL Code | $57.08 \pm 3.75$ | $44.58 \pm 3.54$ | $26.25 \pm 5.18$ | $\mathbf{36.84} \pm 0.44$ |
| **Generator: Qwen3-4B** | | | | |
| Pass@1 | $72.08 \pm 3.96$ | $61.67 \pm 3.98$ | $41.67 \pm 7.13$ | $51.98 \pm 0.58$ |
| Pass@8 | 86.67 | 83.33 | 60.00 | 66.30 |
| Majority@8 | 80.00 | 72.64 | 50.00 | N/A |
| GenSelect 1.7B Prompting | $74.58 \pm 3.96$ | $63.75 \pm 3.75$ | $42.92 \pm 4.52$ | $52.01 \pm 1.09$ |
| GenSelect 4B Prompting | $80.00 \pm 1.78$ | $71.25 \pm 2.48$ | $50.00 \pm 2.52$ | $55.04 \pm 0.88$ |
| GenSelect 8B Prompting | $\underline{82.92} \pm 2.14$ | $72.50 \pm 2.36$ | $53.75 \pm 2.78$ | $\underline{56.11} \pm 1.09$ |
| GenSelect 1.7B RL Math | $\mathbf{80.83} \pm 3.45$ | $\mathbf{\underline{73.33}} \pm 3.09$ | $\mathbf{\underline{55.42}} \pm 1.73$ | $\mathbf{53.36} \pm 0.81$ |
| GenSelect 1.7B RL Code | $77.50 \pm 2.95$ | $67.08 \pm 3.30$ | $47.92 \pm 2.48$ | $52.67 \pm 0.34$ |
| **Generator: Qwen3-8B** | | | | |
| Pass@1 | $76.25 \pm 6.28$ | $69.17 \pm 4.27$ | $40.42 \pm 2.14$ | $55.78 \pm 0.69$ |
| Pass@8 | 90.00 | 86.67 | 56.67 | 66.96 |
| Majority@8 | 83.81 | 76.67 | 50.00 | N/A |
| GenSelect 1.7B Prompting | $78.75 \pm 3.54$ | $72.92 \pm 3.75$ | $46.25 \pm 4.15$ | $55.56 \pm 0.67$ |
| GenSelect 4B Prompting | $83.33 \pm 2.52$ | $77.08 \pm 2.14$ | $52.08 \pm 2.48$ | $57.85 \pm 0.87$ |
| GenSelect 8B Prompting | $\underline{86.25} \pm 2.14$ | $\underline{78.33} \pm 2.52$ | $55.83 \pm 1.54$ | $\underline{58.26} \pm 0.71$ |
| GenSelect 1.7B RL Math | $\mathbf{85.83} \pm 2.36$ | $\mathbf{78.33} \pm 1.78$ | $\mathbf{\underline{56.25}} \pm 1.18$ | $56.33 \pm 0.72$ |
| GenSelect 1.7B RL Code | $79.17 \pm 2.36$ | $74.58 \pm 2.48$ | $49.17 \pm 2.95$ | $\mathbf{56.51} \pm 0.59$ |

Table 1: Accuracy in % when selecting among $N$ candidate solutions generated by different models (Qwen3-1.7B/4B/8B). Pass@8 represents the maximum achievable selection accuracy while RL models compared with equivalent baselines are highlighted in gray. The best method among equivalently sized strategies is **bolded**, and the overall best method is underlined.

## 3.2. Evaluation

For evaluation, we report math results on AIME24, AIME25, and HMMT25, and code results on the Live-CodeBench v6 split (August 2024–May 2025). All evaluations are performed with the NeMo-Skills framework [14]. Prompting baselines use the same selection prompt employed during training (Figure 2). GenSelect results are averaged over 8 runs. Majority@8 corresponds to the accuracy obtained by selecting the most common answer (applicable only for math tasks), and Pass@8 defines the maximum achievable selection accuracy.

## 4. Results

Table 1 reports the accuracy when applying our trained models to select among candidate sets generated by Qwen3-1.7B, Qwen3-4B, and Qwen3-8B.

### 4.1. Results on Qwen3-1.7B Generations

Table 1 (Qwen3-1.7B candidates) shows that reinforcement learning substantially improves GenSelect performance over pass@1, majority voting, and GenSelect prompting with Qwen3-1.7B. In math, the RL-trained 1.7B model surpasses the out-of-the-box per-
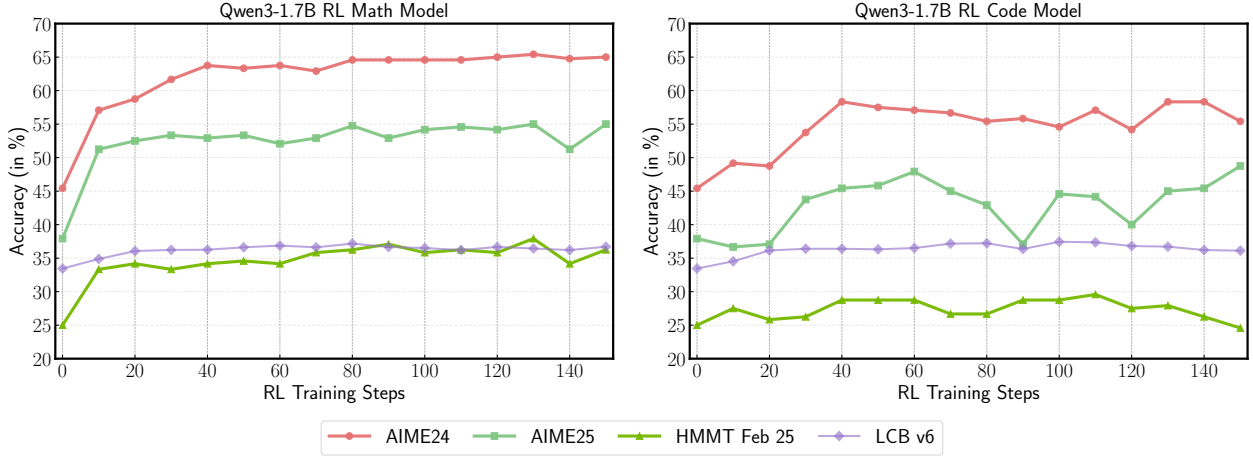
Figure 3: RL training dynamics of `Qwen3-1.7B` for GenSelect on math (left) and code (right) selection tasks. Selection accuracy improves rapidly and stabilizes for math, while training on code shows noisier dynamics and earlier saturation, reflecting domain-dependent differences in selection difficulty.

formance of Qwen3-4B and approaches that of Qwen3-8B across AIME24, AIME25, and HMMT25. In code, the RL-trained 1.7B model consistently outperforms same-size prompting baselines but does not exceed Qwen3-4B or Qwen3-4B GenSelect prompting.

We observe that GenSelect prompting with `Qwen3-1.7B` underperforms its 4B and 8B counterparts noticeably despite identical prompts. While lower values are expected, we hypothesize this is also due to the smaller 1.7B parameter count, which may cause the model to overfit to typical solution-generation tasks and fail to generalize to the GenSelect instructions. With a small amount of RL training, we enable the model to learn solution selection at a level comparable to an 8B-parameter model.

### 4.2. Results on Higher Quality Generations

Table 1 (Qwen3-4B and Qwen3-8B candidates) demonstrates that the gains from reinforcement learning generalize to higher-quality candidate sets. When selecting from Qwen3-4B and Qwen3-8B generations, our RL-trained math models match or exceed the out-of-the-box selection performance of Qwen3-8B on AIME25 and HMMT25, while remaining competitive on AIME24 and LiveCodeBench.

These results are particularly notable because selection over Qwen3-4B and Qwen3-8B generations is off-policy relative to the `Qwen3-1.7B`-generated candidates used during training. This indicates that reinforcement learning enables the model to learn a transferable notion of solution quality rather than overfitting to a specific generation model. Across all generation qualities, GenSelect models trained with RL outperform majority voting and equivalently sized

prompting baselines. This establishes reinforcement learning as an effective strategy for improving parallel thinking and one that generalizes to different generation models.

### 4.3. Math and Coding Model Differences

Performance gains from reinforcement learning are consistently larger for math than for code. While RL-trained code models outperform GenSelect prompting with Qwen3-1.7B, they remain behind 4B and 8B prompting baselines on LiveCodeBench. We also observe asymmetric cross-domain transfer: training on math selection data improves code selection, whereas training on code yields limited improvements on math, consistent with prior findings [3]. Figure 3 shows the accuracy after training of increments of 10 steps steadily improving the math model on both math and coding domains. Alternatively, the coding model demonstrates a similar pattern for LCB but marginal and less stable improvements on math benchmarks.

We attribute the smaller gains in code to two factors. First, correctness signals derived from unit tests are noisier and less exhaustive than symbolic math verification, weakening the reinforcement signal: a buggy program can pass a limited test suite (false positives), and otherwise-correct solutions can fail due to corner cases not covered by the tests or to runtime/formatting issues. Second, code selection requires broader semantic coverage than math—including algorithmic structure, edge-case handling, complexity constraints, and implementation details—making comparative reasoning more difficult than in math, where correctness is often more sharply defined [25].

# 5. Related Work

Test-time compute has emerged as a highly effective method for improving LLM reasoning, and has become a central theme across recent work on scaling reasoning capabilities. Early work on self-consistency [23] demonstrated that sampling multiple reasoning paths and aggregating them via majority voting can substantially improve accuracy. Subsequent surveys and scaling studies [29, 30] have systematized this direction, identifying test-time compute as a promising and complementary alternative to scaling model size.

## 5.1. Reward Models

Reward models have proven effective for rating and selecting candidate solutions, with multiple works highlighting their benefits in mathematical and coding tasks [11, 27]. To further exploit the advantages of test-time compute, recent efforts have investigated generative approaches: verifiers framed as next-token predictors [28], reward models trained on reasoning preferences [12, 8], and models that reason via Chain-of-Thought before scoring. Other directions include pairwise tournament-style judging [10], generating test cases as verifiers [6], and reinforcement-learning verifiers for competitive math proofs [18]. Further refinements leverage critique fine-tuning to enhance scoring [24, 1] and internal model confidence signals to filter reasoning traces [7].

## 5.2. Parallel Reasoning Selection

Parallel reasoning methods extend beyond scoring to structured comparison, aggregation, and refinement over multiple candidate solutions. GenSelect [21] reframes Best-of-$N$ as a generative selection problem, using a reasoning model to explicitly compare and select among candidates. Related approaches explore self-refinement and iterative improvement during training [22], as well as aggregation over parallel outputs using learned reasoning policies [16, 31].

In this work, we focus exclusively on *selection* rather than synthesis. This choice is motivated by the strong empirical performance of GenSelect [21] and findings by Qi et al. [16] showing that synthesized outputs often reduce to copying one of the inputs.

While synthesis strategies are effective, in this work we show that selection can be equally successful. In addition, selection can be substantially easier to train with reinforcement learning: since rewards are defined over selected candidates, correctness labels can be precomputed for the entire candidate set. In contrast, synthesis-based aggregation produces new outputs that typically require re-running verification (or

reward-model evaluation), which can be costly or difficult to integrate in complex RL environments. We extend the work by Toshniwal et al. [21] by proposing training strategies to enhance the GenSelect performance of small reasoning models.

# 6. Conclusion

We study reinforcement learning for GenSelect in small reasoning models on the Best-of-$N$ selection problem in math and code. Training a 1.7B-parameter model directly for generative selection consistently outperforms majority voting and GenSelect prompting, in several cases exceeding the performance of much larger models. The learned selection behavior generalizes beyond the training-time candidate distribution, enabling effective selection among candidates generated by the frozen 1.7B base model at evaluation time. Reinforcement learning yields larger and more reliable gains in math than in code, with stronger cross-domain transfer when training on math. Overall, our results show that training for selection is a useful complementary approach to test-time scaling, enabling small models to exhibit powerful aggregation behavior.

# References

[1] Wasi Uddin Ahmad, Somshubra Majumdar, Aleksander Ficek, Sean Narenthiran, Mehrzad Samadi, Jocelyn Huang, Siddhartha Jain, Vahid Noroozi, and Boris Ginsburg. OpenCodeReasoning-II: A Simple Test Time Scaling Approach via Self-Critique, 2025. URL https://arxiv.org/abs/2507.09075.

[2] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large Language Monkeys: Scaling Inference Compute with Repeated Sampling, 2024. URL https://arxiv.org/abs/2407.21787.

[3] Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee, Peng Xu, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. AceReason-Nemotron: Advancing Math and Code Reasoning through Reinforcement Learning. In *NeurIPS*, 2025.

[4] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training Verifiers to Solve Math Word Problems, 2021. URL https://arxiv.org/abs/2110.14168.

[5] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi 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, Wanjia Zhao, 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, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning, 2025. URL https://arxiv.org/abs/2501.12948.

[6] Aleksander Ficek, Somshubra Majumdar, Vahid Noroozi, and Boris Ginsburg. Scoring Verifiers: Evaluating Synthetic Verification for Code and Reasoning. In *COLM*, 2025.

[7] Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep Think with Confidence, 2025. URL https://arxiv.org/abs/2508.15260.

[8] Jiaxin Guo, Zewen Chi, Li Dong, Qingxiu Dong, Xun Wu, Shaohan Huang, and Furu Wei. Reward Reasoning Model, 2025. URL https://arxiv.org/abs/2505.14674.

[9] Hynek Kydlíček. Math-Verify: Math Verification Library. URL https://github.com/huggingface/math-verify.

[10] Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. PairJudge RM: Perform Best-of-N Sampling with Knockout Tournament, 2025. URL https://arxiv.org/abs/2501.13007.

[11] Zihan Liu, Yang Chen, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. AceMath: Advancing Frontier Math Reasoning with Post-Training and Reward Modeling, 2025. URL https://arxiv.org/abs/2412.15084.

[12] Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. Generative Reward Models, 2024. URL https://arxiv.org/abs/2410.12832.

[13] Ivan Moshkov, Darragh Hanley, Ivan Sorokin, Shubham Toshniwal, Christof Henkel, Benedikt Schifferer, Wei Du, and Igor Gitman. AIMO-2 Winning Solution: Building State-of-the-Art Mathematical Reasoning Models with OpenMathReasoning dataset, 2025. URL https://arxiv.org/abs/2504.16891.

[14] NVIDIA. NeMo-Skills. https://github.com/NVIDIA-NeMo/Skills/, 2025. GitHub repository.

[15] OpenAI. OpenAI o1 System Card, 2024. URL https://arxiv.org/abs/2412.16720.

[16] Jianing Qi, Xi Ye, Hao Tang, Zhigang Zhu, and Eunsol Choi. Learning to Reason Across Parallel Samples for LLM Reasoning, 2025. URL https://arxiv.org/abs/2506.09014.

[17] 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. In *Proceedings of the Twentieth European Conference on Computer Systems*, pages 1279–1297, 2025.

[18] Wenlei Shi and Xing Jin. Heimdall: test-time scaling on the generative verification, 2025. URL https://arxiv.org/abs/2504.10337.

[19] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters. In *ICLR*, 2025.

[20] Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisacanin, Alexan Ayrapetyan, and Igor Gitman. OpenMathInstruct-2: Accelerating AI for Math with Massive Open-Source Instruction Data. In *ICLR*, 2025.

[21] Shubham Toshniwal, Ivan Sorokin, Aleksander Ficek, Ivan Moshkov, and Igor Gitman. GenSelect: A Generative Approach to Best-of-N, 2025. URL https://arxiv.org/abs/2507.17797.

[22] Qibin Wang, Pu Zhao, Shaohan Huang, Fangkai Yang, Lu Wang, Furu Wei, Qingwei Lin, Saravan Rajmohan, and Dongmei Zhang. Learning to Refine: Self-Refinement of Parallel Reasoning in LLMs, 2025. URL https://arxiv.org/abs/2509.00084.

[23] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language Models, 2023. URL https://arxiv.org/abs/2203.11171.

[24] Yubo Wang, Xiang Yue, and Wenhu Chen. Critique Fine-Tuning: Learning to Critique is More Effective than Learning to Imitate, 2025. URL https://arxiv.org/abs/2501.17703.

[25] Zihan Wang, Siyao Liu, Yang Sun, Hongyan Li, and Kai Shen. CodeContests+: High-Quality Test Case Generation for Competitive Programming, 2025. URL https://arxiv.org/abs/2506.05817.

[26] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. DAPO: An Open-Source LLM Reinforcement Learning System at Scale. *arXiv preprint arXiv:2503.14476*, 2025.

[27] Huaye Zeng, Dongfu Jiang, Haozhe Wang, Ping Nie, Xiaotong Chen, and Wenhu Chen. ACE-CODER: Acing Coder RL via Automated Test-Case Synthesis, 2025. URL https://arxiv.org/abs/2502.01718.

[28] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative Verifiers: Reward Modeling as Next-Token Prediction, 2025. URL https://arxiv.org/abs/2408.15240.

[29] Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang, Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo, Yufei Wang, Niklas Muennighoff, Irwin King, Xue Liu, and Chen Ma. A Survey on Test-Time Scaling in Large Language Models: What, How, Where, and How Well?, 2025. URL https://arxiv.org/abs/2503.24235.

[30] Eric Zhao, Pranjal Awasthi, and Sreenivas Gollapudi. Sample, Scrutinize and Scale: Effective Inference-Time Search by Scaling Verification, 2025. URL https://arxiv.org/abs/2502.01839.

[31] Wenting Zhao, Pranjal Aggarwal, Swarnadeep Saha, Asli Celikyilmaz, Jason Weston, and Ilia Kulikov. The Majority is not always right: RL training for solution aggregation, 2025. URL https://arxiv.org/abs/2509.06870.