

# TOWARDS EXECUTION-GROUNDED AUTOMATED AI RESEARCH

Chenglei Si\*, Zitong Yang\*, Yejin Choi, Emmanuel Candès, Diyi Yang, Tatsunori Hashimoto  
Stanford University  
clsi@stanford.edu, zitong@berkeley.edu  
\*Equal contribution

## ABSTRACT

Automated AI research holds great potential to accelerate scientific discovery. However, current LLMs often generate plausible-looking but ineffective ideas. Execution grounding may help, but it is unclear whether automated execution is feasible and whether LLMs can learn from the execution feedback. To investigate these, we first build an automated executor to implement ideas and launch large-scale parallel GPU experiments to verify their effectiveness. We then convert two realistic research problems – LLM pre-training and post-training – into execution environments and demonstrate that our automated executor can implement a large fraction of the ideas sampled from frontier LLMs. We analyze two methods to learn from the execution feedback: evolutionary search and reinforcement learning. Execution-guided evolutionary search is sample-efficient: it finds a method that significantly outperforms the GRPO baseline on post-training, and finds a pre-training recipe that outperforms the nanoGPT baseline on pre-training, all within just ten search epochs. Frontier LLMs often generate meaningful algorithmic ideas during search, but they tend to saturate early and only occasionally exhibit scaling trends. Reinforcement learning from execution reward, on the other hand, suffers from mode collapse. It successfully improves the average reward of the ideator model but not the upper-bound, due to models converging on simple ideas. We also thoroughly analyze the executed ideas and training dynamics.

## 1 INTRODUCTION

We envision automated AI research: LLMs generate research ideas to tackle important research problems, implement the ideas as code, run experiments to verify the effectiveness, and continuously learn from the execution results. If successful, these automated AI researchers could automatically develop and identify effective research ideas in a massive search space, thereby scalably converting compute into scientific discovery. Despite the promise, automated AI research is bottlenecked by the ability of LLMs to generate effective ideas. Si et al. (2025b) and Si et al. (2025a) evaluated the quality of LLM-generated research ideas through large-scale expert review and found that LLM ideas often look convincing but are ineffective after being executed by human researchers.

This highlights the need to ground idea generation in execution. However, obtaining execution results of ideas in an automated and scalable manner is challenging, especially since we are targeting open-ended AI research where any ideas expressible in natural language are within our action space. To tackle this, we design and build a high-throughput automated idea executor that can implement

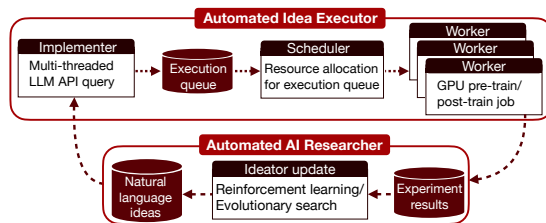


Figure 1: We build an automated idea executor involving Implementer, Scheduler, and Worker. We then use this automated executor as a reward function to teach LLMs to generate more effective ideas through evolutionary search and RL. We only update the ideator in the learning process.

---

hundreds of model-generated ideas and execute them in parallel to obtain the experiment results as execution feedback.

To study the extent to which we can automate realistic LLM research, we chose two GPU-intensive research problems (LLM pre-training and post-training) that are critical for improving the capabilities of LLMs as the research environments for our automated AI researchers. For the first time, we demonstrate that our automated executor can implement a large fraction of LLM-generated ideas on these challenging open-ended research problems, with over 90% execution rates on the pre-training environment with Claude-4.5-Sonnet and Claude-4.5-Opus.

To analyze whether grounding on execution-feedback can improve LLM idea generation, we define objective performance metrics for both environments and analyze the strengths and weaknesses of two popular learning algorithms: evolutionary search and reinforcement learning.

We use our automated executor to guide evolutionary search. Within ten search epochs, this execution-guided search finds a post-training recipe that outperforms the GRPO baseline (69.4% vs 48.0%) on the task of post-training a 1.5B model for math reasoning, and a pre-training recipe that outperforms the nanoGPT baseline (19.7 minutes vs 35.9 minutes) on the task of minimizing the training wall-clock time to reach the target validation loss (Table 1). Our analysis shows that models are often generating algorithmic ideas apart from tuning hyper-parameters, and evolutionary search significantly outperforms best-of-N under the same sampling budget. However, when analyzing the scaling trend, only Claude-4.5-Opus shows a clear scaling curve, while both Claude-4.5-Sonnet and GPT-5 tend to saturate early.

We then use the automated executor as the reward function in an RL loop to finetune Qwen3-30B. We show that RL with execution reward can successfully improve the average reward of the ideator model, similar to typical RL from verifiable rewards. However, RL does not improve the max reward, which is the more important metric for scientific discovery. In fact, we reveal that RL causes the ideator model to converge on a few easy-to-implement ideas, resulting in a collapse in thinking length and idea diversity.

Table 1: Performance of our execution-guided search vs baselines and best human experts.

	Post-Training $\uparrow$	Pre-Training $\downarrow$
Baseline	48.0%	35.9 min
Execution-Guided Search	69.4%	19.7 min
Best Human Expert	68.8%	2.1 min

## 2 AUTOMATED IDEA EXECUTOR

To measure the effectiveness of model-generated ideas, we build an automated executor that takes natural language research ideas as input, generates code implementations, runs the experiments on the backend, and returns the idea’s benchmark performance as the final output.

### 2.1 RESEARCH ENVIRONMENTS FOR IDEATION

Our automated idea executor is grounded in specific research environments, where each environment consists of a research problem, a baseline codebase, a benchmark to measure performance on, fixed training and evaluation data, and evaluation metrics. When constructing the research environments, we aim to select research problems that are open-ended, so that there is ample room for new algorithmic innovations, and at the same time have well-established baselines and benchmarking metrics so that measuring effectiveness is straightforward. In this work, we construct both a pre-training environment and a post-training environment for the automated AI researchers to work on.

**Pre-Training Task: Improving nanoGPT** In the nanoGPT environment, we provide a baseline codebase adapted from the nanoGPT speedrun Jordan et al. (2024) and ask the ideator model to brainstorm possible improvements. The original speedrun task is to minimize the time to pre-train a 124M GPT-2 model Radford et al. (2019) on the FineWeb corpus Penedo et al. (2024) to reach a validation loss of 3.28 on the validation set. We did several modifications to the original speedrun setting. First, we introduce a proxy reward equal to the reciprocal of the validation loss ( $\frac{1}{loss}$ ) when performing the search and RL experiments in later sections of the paper. This way, we can fix the training wall-clock time to be 25 minutes and ask the model to directly optimize the proxy reward under this fixed budget, so that we can avoid different runs having vastly different runtimes. We report

---

the validation loss or the proxy reward metric in most plots, and only measure and report the training time metric for the top solution in order to directly compare it with the human experts’ solutions on the original nanoGPT speedrun leaderboard. Second, to avoid any possible reward hacking, we freeze all evaluation hyper-parameters and implement an inference function that predicts one future token at a time to prevent models from changing the attention mechanism in a way that leaks future tokens (which happened multiple times during our initial development). We use this inference function during the final validation after each training run.

**Post-Training Task: Improving GRPO** In the GRPO environment, the baseline is an implementation of the GRPO algorithm Shao et al. (2024) that finetunes a Qwen2.5-Math-1.5B checkpoint Yang et al. (2024) on the MATH dataset Hendrycks et al. (2021). The ideator model needs to brainstorm post-training algorithms more effective than the baseline. We specify a fixed training wall-clock time budget and use the max accuracy on the MATH validation set during training as the metric. To prevent reward hacking, we keep all validation-related code in a separate file and do not allow the automatic executor to access or modify it.

## 2.2 SYSTEM DESIGN

The automated idea executor can be viewed as a high-level API whose input is a batch of natural language ideas, and the output is the benchmark performance of each idea. There are three core building blocks of this API (Figure 1): *Implementer* – the server that generates the code diff for the idea and applies those changes; *Scheduler* – a middle layer that receives the list of codebases and allocates resources to run experiments; *Worker* – the cluster with GPU available that runs the experiments and uploads the experiment results.

**Implementer** The implementer is hosted on a CPU machine with high IO capacity. First, the user submits a batch of natural language ideas. Then, for each idea, the implementer makes parallelized API calls to the code execution LLM to obtain a `diff` file that can be patched into the corresponding baseline codebase. To optimize for efficiency, we prompt the code execution LLM with both the idea and the baseline codebase to sample 10 code diff files in parallel. For each sample, if the generated diff file cannot be patched into the original codebase, we provide the patch log and ask the model to revise the original generation. We repeat this sequential self-revision for a maximum of 2 times. In the end, we return the first code diff file that can be successfully patched into the baseline codebase. The patched codebase is then submitted to a cloud bucket as a `.zip` file.

**Scheduler** Under a set clock frequency, the scheduler downloads the new codebases from the cloud. If the codebase has not been executed, the scheduler examines the resource requirement of the given research environment and prepares a job configuration to be submitted.

**Worker** Once the scheduler finds available resources, it connects the prepared job configuration with the GPU resource and initializes the worker to run the experiment. If the execution of the experiment is successful, the worker will upload the experiment logs including all performance metrics to another cloud bucket (`wandb`) along with the complete metadata: idea content, code change, execution log, etc. If the execution fails (e.g., due to bugs in code implementation), the worker halts. The user (i.e., the ideator model) can then download the execution results and see the performance of the batch of ideas they submitted with full training logs.

## 3 BENCHMARKING LLM IDEATORS AND EXECUTORS

The prerequisite for an execution-grounded feedback loop is that current LLMs can serve as both ideators and executors, so that we can get meaningful reward signals for the models to learn from. To examine this prerequisite, we benchmark various frontier LLMs as both the ideator and the executor.

### 3.1 END-TO-END IDEATION AND EXECUTION

In the first setting, we sample ideas from an LLM, and use the same LLM as the code execution model to execute its own ideas. We sample and execute 50 ideas from Claude-4.5-Opus, Claude-4.5-Sonnet, and GPT-5, and measure several metrics: (1) completion rate: the percentage of ideas that are successfully executed with a valid (non-zero) experiment result after execution; (2) average performance: the average validation accuracy or loss for all the successfully executed ideas among the

50 samples; (3) best performance: the highest validation accuracy or lowest validation loss among all executed ideas. We present results in the top row of Figure 2. Notably, a large fraction of the sampled ideas can indeed be executed successfully, with Claude-4.5-Opus and Claude-4.5-Sonnet having a significantly higher execution rate than GPT-5. Moreover, the best-of-N performance ( $N = 50$ ) of these models can already beat the original baseline solutions. For example, on the GRPO environment, Claude-4.5-Sonnet gets a max accuracy of 60.4% as compared to the baseline of 48.0%; on nanoGPT, Claude-4.5-Opus gets a lowest loss of 3.237 as compared to the baseline of 3.255.

### 3.2 COMPARING IDEATORS WITH THE SAME EXECUTOR

In the second setting, we fix the executor model to be GPT-5 and use different ideator models to sample ideas. As shown in the bottom row of Figure 2, even when the ideator and executor are different models, the execution rate is still decent (ranging from 42% to 78%), although we do notice that the same ideas from Claude-4.5-Sonnet get a lower execution rate when executed by GPT-5 instead of itself (84% vs 42% on GRPO and 90% vs 78% on nanoGPT). Moreover, frontier open-weight models like Kimi-K2-Thinking Kimi Team (2025) and Qwen3-235B-A22B Yang et al. (2025a) can also get non-trivial completion rates and achieve best-of-N performance that outperforms the baseline solutions in this setting. For example, Qwen3-235B achieves a max accuracy of 50.2% on GRPO and min loss of 3.238 on nanoGPT with  $N = 50$ , both better than the baselines.

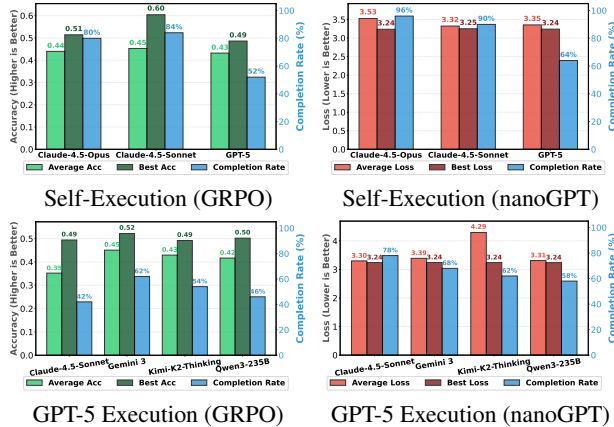


Figure 2: Model performance comparison with self-execution (top row) vs GPT-5 execution (bottom row) on the GRPO and nanoGPT environments. The baseline accuracy for GRPO is 0.480, and the baseline loss for nanoGPT is 3.255.

These benchmarking results demonstrate the feasibility of the automated ideation and execution loop. Next, we build search scaffolds and RL training loops to examine whether models can learn from the execution feedback.

## 4 EXECUTION-GUIDED EVOLUTIONARY SEARCH

We develop an evolutionary search scaffold on top of frontier LLMs to optimize for effective ideas based on execution feedback. We introduce our search method that blends exploration and exploitation, its effectiveness on our two research environments, and various analyses of the generated ideas throughout the evolutionary search process.

### 4.1 SEARCH SCAFFOLD

Our search method is inspired by prior evolutionary search approaches for code optimization, such as AlphaEvolve Novikov et al. (2025). At the first search epoch, we sample a full batch of new ideas. In all subsequent epochs, we split the idea generation into exploitation and exploration subsets. For exploitation, we choose ideas from previous epochs that outperform the baseline and append them to the idea generation prompt to ask the ideator model to generate new variants that combine their strengths. For exploration, we randomly sample ideas from previous epochs to append to the idea generation prompt until reaching the max context length and instruct the ideator model to generate completely new ideas different from them. We start with 50% exploitation and 50% exploration at epoch 1 and gradually anneal the exploration rate and increase the exploitation ratio in later epochs. We use a batch size of 50 for the GRPO environment and a batch size of 80 for the nanoGPT environment.

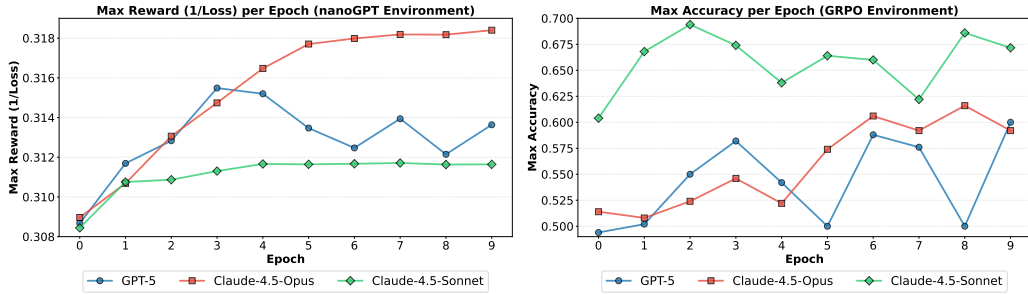


Figure 3: Best performance at each epoch when performing execution-guided search with different models. For the nanoGPT environment (left), we use the reciprocal of the validation loss as the metric; for the GRPO environment (right), we use validation accuracy as the metric. Claude-4.5-Opus exhibits a scaling trend on both environments and achieves the best performance on nanoGPT. Claude-4.5-Sonnet achieves the best performance on GRPO due to effective hyper-parameter tuning, but saturates early.

## 4.2 EXPERIMENT RESULTS

For each environment, we perform execution-guided search with three different models: Claude-4.5-Opus, Claude-4.5-Sonnet, and GPT-5. For each experiment, we use the same model as both the ideator and executor (i.e., self-execution). We plot the progression of the best performance at each search epoch in Figure 3. We summarize several notable trends below.

First, we observe a scaling trend with Claude-4.5-Opus, where searching for more epochs leads to a higher upper bound. In contrast, Claude-4.5-Sonnet and GPT-5 tend to saturate early. Second, all models can find ideas that significantly outperform the baselines. On GRPO, Claude-4.5-Sonnet finds that using vanilla policy gradient with the group-average baseline without importance reweighting or clipping outperforms the standard GRPO objective in this particular experiment setup and exploits this finding in all subsequent search epochs, resulting in the best solution of 69.4% at epoch 2 with precise hyper-parameter tuning. On nanoGPT, Claude-4.5-Opus achieves the min validation loss of 3.1407 at epoch 9 by combining various architectural modifications, hyper-parameter tuning, and applying exponential moving average of intermediate checkpoints during validation (see Appendix A.4 for the full idea). We run this top solution on 8 H100s to follow the same setup as the nanoGPT speedrun, and it reaches the 3.28 target validation loss in 19.7 minutes, a significant speedup as compared to the baseline codebase, which takes 35.9 minutes of training time to reach the same target validation loss.

To better contextualize these solutions optimized by the model, we also compare the top performance of our execution-guided search to human experts. For the GRPO environment, we compare with the leaderboard of a graduate-level LLM class, which hosted the same environment as an assignment for all students to optimize the validation accuracy under the same training time budget. The best student solution achieved an accuracy of 68.8%, lower than Claude-4.5-Sonnet’s top solution using our execution-guided search. For the nanoGPT environment, we directly compare with the nanoGPT speedrun leaderboard,<sup>1</sup> where the top human solution as of December 2025 can reach the target validation loss under 2.1 minutes, indicating significant headroom for further model capability and search method improvement on this environment.

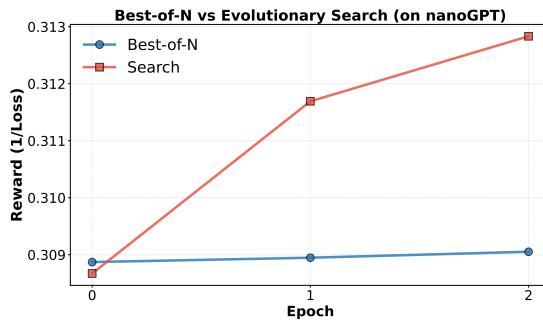


Figure 4: Comparison between best-of-N and our execution-guided search under the same sampling budget.

<sup>1</sup><https://github.com/KellerJordan/modded-nanogpt>

Table 2: Breakdown of hyper-parameter tuning vs algorithmic ideas throughout the entire execution-guided search. We report the percentage of each type among all generated ideas of each model ( $N = 500$  ideas on GRPO and  $N = 800$  ideas on nanoGPT). Bold numbers every row indicate the best performance by each model. All models generate a substantial amount of algorithmic ideas apart from hyper-parameter changes, while Claude-4.5-Sonnet generates significantly more hyper-parameter ideas than other models.

Model name	Hyper-parameter			Algorithmic		
	Percentage	Average	Best	Percentage	Average	Best
<i>GRPO environment (accuracy<math>\uparrow</math>)</i>						
GPT-5	5.0%	45.0%	50.2%	95.0%	44.5%	<b>60.0%</b>
Claude-4.5-Sonnet	41.1%	48.4%	<b>69.4%</b>	58.9%	45.0%	67.4%
Claude-4.5-Opus	3.7%	44.4%	50.4%	96.3%	46.5%	<b>61.6%</b>
<i>nanoGPT environment (loss<math>\downarrow</math>)</i>						
GPT-5	15.4%	3.254	3.195	84.6%	3.894	<b>3.170</b>
Claude-4.5-Sonnet	31.3%	3.251	<b>3.208</b>	68.7%	3.679	<b>3.208</b>
Claude-4.5-Opus	8.7%	3.329	3.147	91.3%	3.419	<b>3.141</b>

### 4.3 COMPARISON WITH BEST-OF-N

To demonstrate the effectiveness of our search scaffold, we compare our execution-guided search with the best-of-N baseline with the same sampling budget on the nanoGPT environment. Since the batch size for our search is 80, we compare the first 3 epochs of the execution-guided search using the GPT-5 backbone with the best-of-N results of GPT-5 with  $N \in \{80, 160, 240\}$ . As shown in Figure 4, search and best-of-N start from similar performance at epoch 0 (they are not exactly the same due to variances from sampling), but evolutionary search significantly outperforms best-of-N from epoch 1 onward, demonstrating that the model is effectively leveraging trajectories from previous epochs to generate more effective ideas in future epochs.

### 4.4 ANALYSIS OF GENERATED IDEAS

To quantitatively understand the types of ideas that models generate during the execution-guided search, we perform a stratified analysis by classifying all generated ideas into either hyper-parameter tuning (including any ideas that can be implemented via changing existing configs) or algorithmic (including all ideas that involve implementing new changes not originally supported by the baseline codebase) by using an LLM-judge. Based on Table 2, all three models generate a substantial amount of algorithmic ideas apart from hyper-parameter tuning. Interestingly, different models exhibit different patterns, where Claude-4.5-Sonnet generates significantly more hyper-parameter ideas than both Claude-4.5-Opus and GPT-5. Moreover, the most effective ideas come from algorithmic ideas in most cases, except when using Claude-4.5-Sonnet. To complement the quantitative analysis, we provide several executed ideas in Table 3, as well as Appendix A.4 and Appendix A.5.

## 5 REINFORCEMENT LEARNING FROM EXECUTION REWARD

Different from evolutionary search, reinforcement learning is an alternative learning algorithm that shapes model behaviors through gradient updates. For the first time, we explore whether we can leverage the automated executor as a reward function to directly finetune LLMs to generate more effective ideas via RL. We detail our implementation, experiment setup, and analysis of the training dynamics.

### 5.1 REWARD DESIGN AND EXPERIMENT SETUP

We use Qwen3-30B-A3B Yang et al. (2025a) as the base model and finetune it using the standard GRPO algorithm Shao et al. (2024), motivated by its consistent empirical success on other verified domains. Our prompt batch size is one since we only have one prompt for each environment. In the

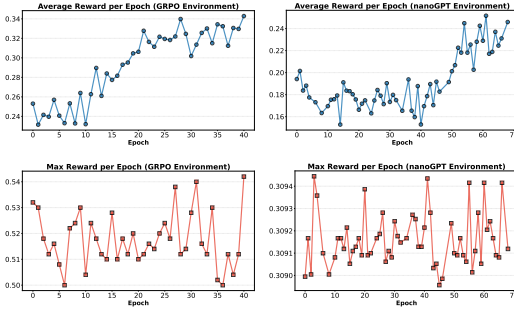


Figure 5: Training curves of RL from execution reward. We plot the average reward per epoch in the upper row, and the max reward per epoch in the lower row. For GRPO, reward is accuracy; for nanoGPT, reward is reciprocal loss. The average reward increases, but not the max reward.

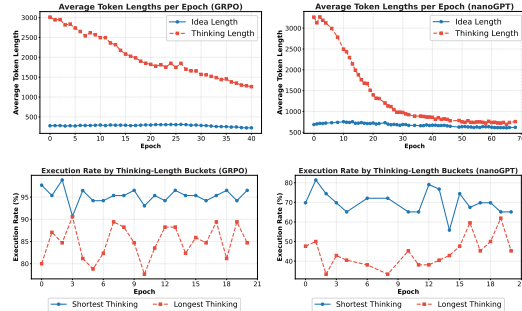


Figure 6: Upper Row: Average length of the thinking trace and the idea per training epoch. Lower Row: Execution rate of ideas with the longest versus shortest thinking traces. Ideas with longer thinking have lower execution rates; correspondingly, thinking lengths decrease in RL.

prompt, we provide the baseline codebase and ask the model to generate new ideas to improve the baseline.

We use large group sizes to stabilize training. For the post-training environment, we use a group size of 256; for the pre-training environment, we use a group size of 128. Since each idea on the GRPO environment runs on one single GPU and each idea on the nanoGPT environment runs on 8 GPUs, these group sizes correspond to parallel execution on 256 GPUs (for GRPO) or 1024 GPUs (for nanoGPT) to obtain the execution reward on each batch of rollout ideas. Each rollout is a thinking trace followed by the natural language idea. We set a max output length of 8192 tokens for rollout sampling and only feed the extracted ideas to the automated executor without the preceding thinking trace.

For the post-training environment, we use the validation set accuracy of each rollout idea after execution as the reward. For ideas without a valid accuracy (i.e., when the execution failed due to code generation errors), we assign a reward of 0. For the pre-training environment, we use the reciprocal of the validation loss as the reward ( $\frac{1}{loss}$ ) and assign a reward of 0 to ideas with failed execution. Our experiments are based on Tinker Thinking Machines Lab (2025).

## 5.2 EXPERIMENT RESULTS

**Positive Training Curves for Average Reward** We plot the average reward of all rollouts of each training epoch in the upper row of Figure 5. For the first time, we demonstrate that the average performance of the generated ideas can increase after sufficient training epochs for open-ended research environments. For instance, the average accuracy on the GRPO environment increases from 0.253 at the beginning to 0.343 after 40 training epochs (top left plot of Figure 5); and the average reward on the nanoGPT environment increases from 0.194 at the beginning to 0.246 after 68 epochs (top right plot of Figure 5), corresponding to a decrease in the average validation loss from 5.150 to 4.066.

**The Case of Max Reward** Despite successfully reproducing the positive training curves observed in other domains, we argue that there is a distinction between idea generation and other verifiable domains. For advancing scientific discovery, we often care about the upper-bound of idea generation, rather than the average quality. In our particular case, we care more about having one breakthrough idea that dominates the baselines, rather than having many safe ideas with a high average. Thus, we plot the max reward of all rollouts at each training epoch in the lower row of Figure 5. The trend here is very different – **the max reward is fluctuating throughout RL training without a clear upward trend**. This reveals the crucial limitation of the standard GRPO algorithm for improving idea generation. In the next subsection, we analyze why RL from execution reward improves the average reward but not the max.

### 5.3 ANALYSIS OF TRAINING DYNAMICS

**Thinking Length** We first plot how the lengths of the thinking traces evolve over RL training in the upper row of Figure 6. In both environments, the thinking traces rapidly decrease in length while the idea lengths stay roughly constant. This is the opposite of the thinking emergence trend from prior RLVR work, such as DeepSeek-R1 DeepSeek-AI et al. (2025). To further understand why the thinking length decreases, we investigate the correlation between the idea execution rate and the thinking trace length. In the bottom row of Figure 6, for the first 20 epochs of the RL training, we sort all ideas in each epoch by their thinking trace lengths and plot the average execution rate of the top-30% longest thinking ideas (red line) and the bottom-30% shortest thinking ideas (blue line). We see a clear trend where ideas with longer thinking consistently have a lower execution rate. We thus hypothesize that longer thinking correlates with more complex ideas with lower execution rates, leading the model to prefer shorter thinking instead in order to maximize the reward.

**Diversity Collapse** Upon manual investigation of all the rollouts being sampled throughout the RL training, we also observed a diversity collapse. Specifically, the model learned to converge on a few simple ideas that can consistently get a positive reward. For example, in the nanoGPT environment, the model learned to converge towards two common ideas: (1) replacing RMSNorm with LayerNorm; and (2) performing exponential moving average (EMA) over intermediate model checkpoints. As shown in Figure 7, out of a batch of 128 sampled ideas per epoch, 51 ideas sampled from Qwen3-30B at epoch 0 are one of the two common ideas above. Towards the end of the RL training, 119 out of 128 sampled ideas at epoch 68 are one of the two common ideas, indicating a severe diversity collapse.

The above analyses reveal that RL causes the models to converge on a few simple-to-implement ideas, accompanied by shrinking thinking lengths. These lead to an increase in the average reward, but do not push the upper-bound due to the lack of exploration. This phenomenon is analogous to mode-collapse observations on other verifiable domains, where the pass@k performance stagnates or even decreases after RL Yue et al. (2025); Wu et al. (2025). Avoiding such convergence and collapse is an open problem and likely requires new algorithmic interventions beyond standard GRPO, which is beyond the scope of this work. However, we do share several preliminary attempts, including: sampling and appending previous epochs’ trajectories into the current epoch’s prompt for rollout sampling, adding a weighted length reward in the total reward, and adding a weighted similarity penalty in the total reward. We did not observe clear gains in the initial epochs and thus early-stopped them, but we document all these results in Appendix A.3 to inform future work.

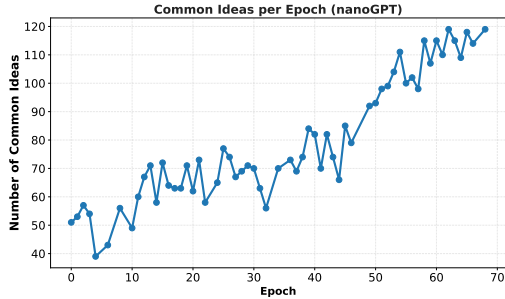


Figure 7: We plot how many ideas in each epoch are about either replacing RMSNorm with LayerNorm or doing EMA. The model converged on these two ideas after RL training.

## 6 CONCLUSION

In this work, we built a large-scale parallel executor for automatically executing model-generated ideas to verify their effectiveness on open-ended LLM research problems, including both pre-training and post-training. Using this executor as a reward function, we analyzed the effectiveness of execution-guided evolutionary search, where frontier LLMs equipped with our search scaffold can significantly outperform the baseline solutions. We also analyzed the limitations of reinforcement learning with execution rewards, where models converge on simple ideas to improve the average reward but lose diversity and do not improve the upper-bound. Our empirical results demonstrate the feasibility and potential of the automated execution feedback loop and also point out the remaining limitations for future improvement.

---

## REFERENCES

- Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi, Amanpreet Singh, Joseph Chee Chang, Kyle Lo, Luca Soldaini, Sergey Feldman, Mike D'Arcy, David Wadden, Matt Latzke, Minyang Tian, Pan Ji, Shengyan Liu, Hao Tong, Bohao Wu, Yanyu Xiong, Luke S. Zettlemoyer, Graham Neubig, Dan Weld, Doug Downey, Wen tau Yih, Pang Wei Koh, and Hanna Hajishirzi. OpenScholar: Synthesizing Scientific Literature with Retrieval-augmented LMs. *ArXiv*, abs/2411.14199, 2024.
- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Mkadry. MLE-bench: Evaluating Machine Learning Agents on Machine Learning Engineering. In *ICLR*, 2025.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Symbolic Discovery of Optimization Algorithms. In *NeurIPS*, 2023.
- Junyan Cheng, Peter Clark, and Kyle Richardson. Language Modeling by Language Models. In *NeurIPS*, 2025.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Jun-Mei Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiaoling Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bing-Li Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dong-Li 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, Jiong 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, M. 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, Ruiqi Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shao-Kang Wu, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wen-Xia Yu, Wentao Zhang, Wangding 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, Xiangyu Jin, Xi-Cheng 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, Yi Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yu-Jing Zou, Yujia He, Yunfan Xiong, Yu-Wei Luo, Yu mei You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yao Li, Yi Zheng, Yuchen Zhu, Yunxiang Ma, Ying Tang, Yukun Zha, Yuting Yan, Zehui Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhen guo Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zi-An 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. *Nature*, 2025.
- Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Taco Cohen, and Gabriele Synnaeve. RLEF: Grounding Code LLMs in Execution Feedback with Reinforcement Learning. In *ICML*, 2025.
- Shashwat Goel, Rishi Hazra, Dulhan Hansaja Jayalath, Timon Willi, Parag Jain, William F. Shen, Ilias Leontiadis, Francesco Barbieri, Yoram Bachrach, Jonas Geiping, and Chenxi Whitehouse. Training AI Co-Scientists Using Rubric Rewards. *ArXiv*, abs/2512.23707, 2025.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring Mathematical Problem Solving With the MATH Dataset. In *NeurIPS*, 2021.
- Shengran Hu, Cong Lu, and Jeff Clune. Automated Design of Agentic Systems. In *ICLR*, 2025.

- 
- Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger, Sang T. Truong, Weixin Liang, Fan-Yun Sun, and Nick Haber. ResearchCodeBench: Benchmarking LLMs on Implementing Novel Machine Learning Research Code. In *NeurIPS*, 2025.
- Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. AIDE: AI-Driven Exploration in the Space of Code. *ArXiv*, abs/2502.13138, 2025.
- Keller Jordan, Jeremy Bernstein, Brendan Rappazzo, @fernbear.bsky.social, Boza Vlado, You Jiacheng, Franz Cesista, Braden Koszarsky, and @Grad62304977. modded-nanogpt: Speedrunning the nanogpt baseline, 2024. URL <https://github.com/KellerJordan/modded-nanogpt>.
- Kimi Team. Kimi K2: Open Agentic Intelligence. *ArXiv*, abs/2507.20534, 2025.
- Jakub L'ala, Odhran O'Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G. Rodrigues, and Andrew D. White. PaperQA: Retrieval-Augmented Generative Agent for Scientific Research. *ArXiv*, abs/2312.07559, 2023.
- Boaz Lavon, Shahar Katz, and Lior Wolf. Execution Guided Line-by-Line Code Generation. In *NeurIPS*, 2025.
- Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang, Daisy Ding, Xinyu Yang, Kailas Vodrahalli, Siyu He, Daniel Scott Smith, Yian Yin, Daniel A. McFarland, and James Zou. Can large language models provide useful feedback on research papers? a large-scale empirical analysis. *NEJM AI*, 2024.
- Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical Representations for Efficient Architecture Search. In *ICLR*, 2018.
- Yixiu Liu, Yang Nan, Weixian Xu, Xiangkun Hu, Lyumanshan Ye, Zhen Qin, and Pengfei Liu. AlphaGo Moment for Model Architecture Discovery. *ArXiv*, abs/2507.18074, 2025.
- Chris Lu, Samuel Holt, Claudio Fanconi, Alex J. Chan, Jakob Nicolaus Foerster, Mihaela van der Schaar, and Robert Tjarko Lange. Discovering Preference Optimization Algorithms with and for Large Language Models. In *NeurIPS*, 2024a.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Nicolaus Foerster, Jeff Clune, and David Ha. The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery. *ArXiv*, abs/2408.06292, 2024b.
- Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi, Abhijeetsingh Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. Discovery-Bench: Towards Data-Driven Discovery with Large Language Models. In *ICLR*, 2025.
- Ludovico Mitchener, Angela Yiu, Benjamin Chang, Mathieu Bourdenx, Tyler Nadolski, Arvis Sulovari, Eric C. Landsness, Dániel L. Barabási, Siddharth Narayanan, Nicky Evans, Shriya Reddy, Martha S. Foiani, Aizad Kamal, Leah P. Shriver, Fang Cao, Asmamaw T. Wassie, Jon M. Laurent, Edwin Melville-Green, Mayk Caldas Ramos, Albert Bou, Kaleigh F. Roberts, Sladjana Zagorac, Timothy C. Orr, Miranda E. Orr, Kevin J. Zwezdaryk, Ali E. Ghareeb, Laurie McCoy, Bruna Gomes, Euan A Ashley, Karen E. Duff, Tonio Buonassisi, Tom Rainforth, Randall J. Bateman, Michael Skarlinski, Samuel G. Rodrigues, Michaela M. Hinks, and Andrew D. White. Kosmos: An AI Scientist for Autonomous Discovery. *ArXiv*, abs/2511.02824, 2025.
- Deepak Nathani, Lovish Madaan, Nicholas Roberts, Niko Ilay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, Dieuwke Hupkes, Ricardo Silveira Cabral, Tatiana Shavrina, Jakob Foerster, Yoram Bachrach, William Yang Wang, and Roberta Raileanu. MLGym: A New Framework and Benchmark for Advancing AI Research Agents. In *COLM*, 2025.
- Alexander Novikov, Ngán V-u, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav M. Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Pushmeet Kohli, Matej Balog, and Google Deepmind. Alphaevolve: A coding agent for scientific and algorithmic discovery. *ArXiv*, abs/2506.13131, 2025.

- 
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben Allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro von Werra, and Thomas Wolf. The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale. *ArXiv*, abs/2406.17557, 2024.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. 2019.
- Esteban Real, Chen Liang, David R. So, and Quoc V. Le. AutoML-Zero: Evolving Machine Learning Algorithms From Scratch. In *ICML*, 2020.
- Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent Laboratory: Using LLM Agents as Research Assistants. In *Findings of EMNLP*, 2025.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and Daya Guo. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *ArXiv*, abs/2402.03300, 2024.
- Chenglei Si, Tatsunori Hashimoto, and Diyi Yang. The Ideation-Execution Gap: Execution Outcomes of LLM-Generated versus Human Research Ideas. *ArXiv*, abs/2506.20803, 2025a.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs Generate Novel Research Ideas? A Large-Scale Human Study with 100+ NLP Researchers. In *ICLR*, 2025b.
- David R. So, Chen Liang, and Quoc V. Le. The Evolved Transformer. In *ICML*, 2019.
- Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, Johannes Heidecke, Amelia Glaese, and Tejal Patwardhan. PaperBench: Evaluating AI’s Ability to Replicate AI Research. In *ICML*, 2025.
- Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang. AI-Researcher: Autonomous Scientific Innovation. In *NeurIPS*, 2025.
- Thinking Machines Lab. Announcing Tinker, 2025.
- Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yaohui Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Min Zhu, Kilian Adriano Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, E. A. Huerta, and Hao Peng. SciCode: A Research Coding Benchmark Curated by Scientists. *ArXiv*, abs/2407.13168, 2024.
- Edan Toledo, Karen Hambardzumyan, Martin Josifoski, Rishi Hazra, Nicolas Mario Baldwin, Alexis Audran-Reiss, Michael Kuchnik, Despoina Magka, Minqi Jiang, Alisia Maria Lupidi, Andrei Lupu, Roberta Raileanu, Kelvin Niu, Tatiana Shavrina, Jean-Christophe Gagnon-Audet, Michael Shvartsman, Shagun Sodhani, Alexander H. Miller, Abhishek Charnalia, Derek Dunfield, Carole-Jean Wu, Pontus Stenertorp, Nicola Cancedda, Jakob Nicolaus Foerster, and Yoram Bachrach. AI Research Agents for Machine Learning: Search, Exploration, and Generalization in MLE-bench. *ArXiv*, abs/2507.02554, 2025.
- Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. SciMON: Scientific Inspiration Machines Optimized for Novelty. In *ACL*, 2024.
- Hjalmar Wijk, Tao Roa Lin, Joel Becker, Sami Jawhar, Neev Parikh, Thomas Broadley, Lawrence Chan, Michael Chen, Joshua Clymer, Jai Dhyani, Elena Elicheva, Katharyn Garcia, Brian Goodrich, Nikola Jurkovic, Megan Kinniment, Aron Lajko, Seraphina Nix, Lucas Jun Koba Sato, William Saunders, Maksym Taran, Ben West, and Elizabeth Barnes. RE-Bench: Evaluating frontier AI R&D capabilities of language model agents against human experts. *ArXiv*, abs/2411.15114, 2024.
- Fang Wu, Weihao Xuan, Ximing Lu, Zaid Harchaoui, and Yejin Choi. The Invisible Leash: Why RLVR May Not Escape Its Origin. *ArXiv*, abs/2507.14843, 2025.

- 
- Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Nicolaus Foerster, Jeff Clune, and David Ha. The AI Scientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search. *ArXiv*, abs/2504.08066, 2025.
- 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. Qwen2.5-Math Technical Report: Toward Mathematical Expert Model via Self-Improvement. *ArXiv*, abs/2409.12122, 2024.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Jingren Zhou, Junyan Lin, Kai Dang, Keqin Bao, Ke-Pei Yang, Le Yu, Li-Chun Deng, Mei Li, Min Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shi-Qiang Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 Technical Report. *ArXiv*, abs/2505.09388, 2025a.
- Sherry Yang, Joy He-Yueya, and Percy Liang. Reinforcement Learning for Machine Learning Engineering Agents. *ArXiv*, abs/2509.01684, 2025b.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does Reinforcement Learning Really Incentivize Reasoning Capacity in LLMs Beyond the Base Model? In *NeurIPS*, 2025.
- Michael R. Zhang, Nishkrit Desai, Juhan Bae, Jonathan Lorraine, and Jimmy Ba. Using Large Language Models for Hyperparameter Optimization. *ArXiv*, abs/2312.04528, 2023.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xuelling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang Yue. OpenCodeInterpreter: Integrating Code Generation with Execution and Refinement. In *Findings of ACL*, 2024.
- Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. DeepReview: Improving LLM-based Paper Review with Human-like Deep Thinking Process. In *ACL*, 2025.
- Barret Zoph and Quoc V. Le. Neural Architecture Search with Reinforcement Learning. In *ICLR*, 2017.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning Transferable Architectures for Scalable Image Recognition. In *CVPR*, 2017.

---

## A APPENDIX

### A.1 RELATED WORK

**AutoML** Our work has deep connections to the AutoML literature. For example, the Neural Architecture Search (NAS) line of work typically defines a constrained set of neural network operators and optimizes for architectures based on validation set performance through reinforcement learning Zoph & Le (2017); Zoph et al. (2017) or search Liu et al. (2018); So et al. (2019). More recent works explored directly using LLMs to propose architecture variants and implement them for validation Liu et al. (2025); Cheng et al. (2025). Beyond architectures, similar automatic optimizations have been applied to improve hyperparameter tuning Zhang et al. (2023), discover machine learning algorithms Real et al. (2020), improve post-training objectives Lu et al. (2024a), discover better neural network optimizers Chen et al. (2023), and design agent scaffolds Hu et al. (2025). Different from this line of work, we tackle automated AI research in a fully open-ended setting without any constraint on the type of ideas. Moreover, our goal is to improve the effectiveness of idea generation, where natural language ideas represent a higher level of abstraction than specific architecture variants or code optimizations.

**LLM-based Research Agents** Recent works have been building end-to-end research agents Lu et al. (2024b); Yamada et al. (2025); Tang et al. (2025); Schmidgall et al. (2025) that use LLMs to generate ideas and implement them through carefully designed scaffolds. They address open-ended AI research as we do, but do not study how to learn from execution feedback. On the other hand, on more grounded benchmarks with clear performance metrics such as MLE-Bench Chan et al. (2025), RE-Bench Wijk et al. (2024), and ML-Gym Nathani et al. (2025), various works have explored how to learn from execution feedback through search Toledo et al. (2025); Jiang et al. (2025) or RL Yang et al. (2025b) to optimize performance on these targeted ML engineering tasks. While we also study algorithms for learning from execution feedback, we tackle open-ended research problems like pre-training and post-training rather than ML engineering tasks that heavily depend on feature engineering and hyper-parameter tuning.

**AI for Research** Apart from fully end-to-end automated AI research, many works have studied how to use LLMs for specific components of the scientific research pipeline, such as literature review Asai et al. (2024); L’ala et al. (2023), idea generation Si et al. (2025b); Wang et al. (2024), data analysis Majumder et al. (2025); Mitchener et al. (2025), experiment plan generation Goel et al. (2025), research code execution Starace et al. (2025); Hua et al. (2025); Tian et al. (2024), and paper reviewing Liang et al. (2024); Zhu et al. (2025). Our work focuses on automated idea execution and learning from the execution feedback. We consider our work complementary to many of the above works that improve other aspects of the scientific research pipeline.

**Execution Grounding for Code** The idea of learning from execution feedback has been explored in the code generation domain. For example, Zheng et al. (2024) curate data and train models to refine code from either human or execution feedback; Gehring et al. (2025) use end-to-end RL training to teach models to improve code based on execution feedback; Lavon et al. (2025) directly guide code generation with execution signals during inference time. In contrast, our work explores execution grounding for the application of idea generation, where the verification is more complicated and expensive.

---

## A.2 DISCUSSION

Despite encouraging initial signals, there are still many limitations to our current set of experiments, which would be great directions for future improvement.

First, our current procedure does not test the generalizability of the generated ideas. It is possible that the best-performing ideas at the small scales may not transfer to gains at a larger scale or on other datasets. Future work should explore methods that explicitly test such generalizability and scalability, and even incorporate them as part of the optimization objectives.

Second, we have shown that RL with execution reward in our current setup can only improve the average reward but not the upper bound. There are many possible reasons, such as a lack of diversity in the base model and missing exploration incentives in the current RL objective. Future work should explore remedies and better learning algorithms for LLMs to more efficiently learn from the execution feedback. For instance, future works could explore how to exploit richer learning signals from the execution trajectories beyond just scalar rewards.

Third, our current experiment scope is bounded by the capability of the execution agent. There exist many promising model-generated ideas that could not be successfully executed by the execution agent (e.g., see the end of Appendix A.4), leading to noise in the reward signal. Future work could develop more capable execution agents and extend our setup to even more complex research problems. For instance, instead of directly prompting an LLM for code diff, future work can implement more capable coding agents with access to external tools and the ability to install new libraries in the execution environments.

Last but not least, we only explored effectiveness as the training reward in this work. There are many other, more subjective alternative metrics that could complement the effectiveness reward, such as the idea novelty and interestingness. Future work could explore how to computationally measure them and incorporate them as part of the training objective to discover more insightful ideas.

### A.3 OTHER RL ATTEMPTS

We present several attempts to improve our RL from the execution reward recipe.

**Attempt 1: Dynamic Prompt** At each epoch (except the first epoch), we randomly sample different executed idea trajectories from the previous epoch and append them to the idea sampling prompt when sampling new rollouts. This merges in-context learning with RL and adds diversity to the idea sampling process. We present the experiment results on the GRPO environment in Figure 8. We did not see significant improvement in early epochs and thus early stopped.

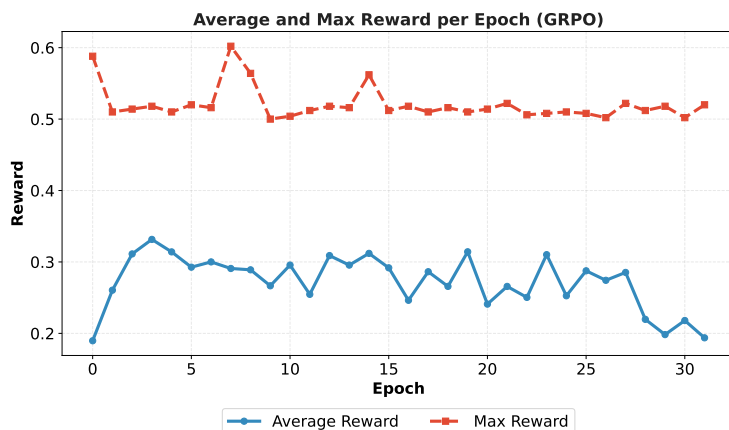


Figure 8: RL with dynamic prompt by appending trajectories from previous epochs. Results are on the GRPO environment.

**Attempt 2: Length Reward** Since we noted a rapid thinking length decrease in our main RL experiment, we tried a simple fix by adding a weighted length reward that counts the number of tokens in the entire rollout sequence, including the thinking trace and the idea. We cap the length reward to a maximum of 0.3 to avoid it dominating the accuracy reward. We present the experiment results on the GRPO environment in Figure 9. As shown on the right panel, the thinking length no longer decreases after adding the length reward to the total reward, but the total training reward isn't going up as shown on the left panel.

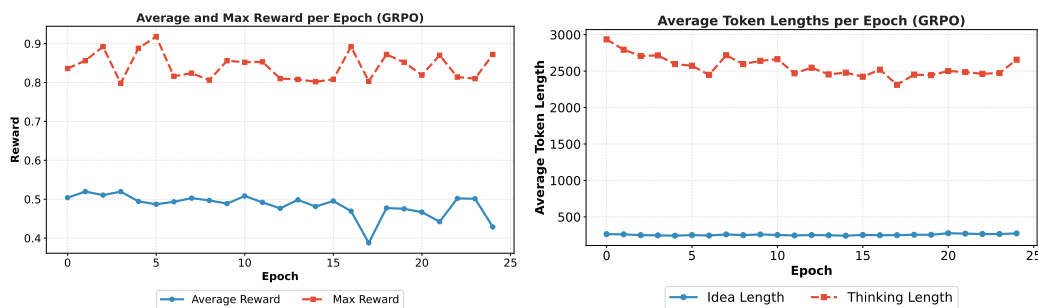


Figure 9: Average and max reward throughout RL training when adding the length reward (left), as well as the progression of the thinking and idea lengths (right).

**Attempt 3: Diversity Reward** We also tried adding a diversity reward in addition to the effectiveness reward. Specifically, when computing the reward for each rollout, we compute its token-level Jaccard similarity with ideas from the previous epoch and add the negative similarity as a penalty to the total reward to discourage repeating ideas that have already been generated before. In fact, this is similar to one of the post-training ideas proposed by Claude-4.5-Sonnet (see example 4 in Appendix A.5). We show the training curves on the GRPO environment in Figure 10. The model maintains a consistent idea similarity (right plot), suggesting sustained exploration. The effectiveness reward is generally showing an upward trend (left plot), but not markedly better than the main RL run with just the effectiveness reward (first sub-plot in Figure 5).

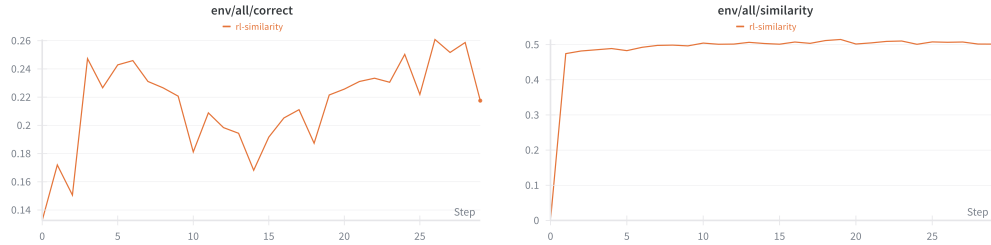


Figure 10: Effectiveness reward (left) and average idea similarity to previous epoch (right) when doing RL with diversity reward.

#### A.4 ADDITIONAL IDEA EXAMPLES

We provide several example ideas on the GRPO environment, including ideas with failed code execution. In most cases, code execution errors happen when the idea involves complicated changes or installing and importing external packages not supported in our execution environment. Future work should explore improvement to the execution agent so that more complicated types of ideas (e.g., training additional auxiliary models or system-level optimizations) can be implemented correctly.

Claude-4.5-Opus on GRPO	Claude-4.5-Sonnet on GRPO
<p><b>Residual Ratio Learning with Momentum Bounds:</b> Instead of directly using the (importance sampling) ratio, decompose it into a “base” component (EMA of batch mean ratios) and a “residual” component (ratio – base). Apply sigmoid bounding only to the residual, allowing the base to capture systematic policy drift while controlling deviations from it. Combined with momentum clip adaptation. Formula: <math>\text{residual} = \text{ratio} - \text{ema\_batch\_ratio}</math>, <math>\text{bounded\_residual} = \text{sigmoid\_bound}(\text{residual}, \text{deviation})</math>, <math>\text{effective\_ratio} = 1.0 + \text{bounded\_residual}</math>.</p> <p><b>Validation Accuracy: 61.6</b></p> <p><b>Advantage Rank Difference Weighting:</b> Instead of using absolute advantage magnitude, weight samples by how far their rank differs from their expected rank under uniform distribution. Samples that significantly outperform or underperform their “expected” position get higher weights. This is distribution-free and robust to outliers. Formula: <math>\text{expected\_rank} = (N-1)/2</math>, <math>\text{rank\_diff} =  \text{actual\_rank} - \text{expected\_rank} </math>, <math>\text{weight} = 0.5 + 0.5 * \text{rank\_diff}</math>.</p> <p><b>Validation Accuracy: 59.2</b></p>	<p><b>Dynamic Mathematical Problem Difficulty Balancing with Performance Feedback:</b> Implement intelligent difficulty balancing that dynamically adjusts the mix of problem difficulties based on recent performance trends. When performance is strong, increase difficulty proportion; when struggling, provide more foundational problems. Combine with the proven hyper-parameters for optimal learning progression.</p> <p><b>Validation Accuracy: 64.0</b></p> <p>Implement token-level reward attribution by using attention weights to identify which input tokens contributed most to correct answers, then amplifying the gradient updates for those tokens during policy gradient training.</p> <p><b>Validation Accuracy: 45.2</b></p> <p>Create mathematical working memory simulation by maintaining a context buffer of mathematical facts, definitions, and intermediate results during problem solving. This buffer gets updated as the model works through problems and provides additional context for subsequent mathematical steps, simulating how humans maintain mathematical working memory during complex calculations.</p> <p><b>Validation Accuracy: 58.0</b></p>

Table 3: Examples of successfully executed ideas on the GRPO environment, along with their accuracy on the MATH validation set. The baseline accuracy is 48.0% on this environment.

Successful Execution	Failed Execution
<p><b>[Experiment]</b> Sequence Position Weighted Trust Region: Apply tighter sigmoid bounds to earlier tokens in the sequence (where errors compound) and looser bounds to later tokens. Weight: <math>\text{position\_weight} = 1 - 0.3 * (\text{position} / \text{seq\_len})</math>, <math>\text{effective\_deviation} = 0.25 + 0.2 * \text{position\_weight}</math>. This accounts for the sequential nature of autoregressive generation.</p> <p><b>[Code Changes]</b> Modify <code>grpo.py</code>: Initialize <code>current_cliprange = 0.2</code>, <code>ema_clip_fraction = 0.15</code>. Standard momentum clip updates. Modify <code>compute_grpo_clip_loss</code> in <code>grpo_utils.py</code>: After computing ratio on line 91 (shape: <code>batch_size x seq_len</code>): <code>batch_size, seq_len = ratio.shape, positions = torch.arange(seq_len, device=ratio.device).float().unsqueeze(0).expand(batch_size, -1)</code>, <code>position_weight = 1.0 - 0.3 * (positions / (seq_len - 1 + 1e-6))</code>, <code>effective_deviation = 0.25 + 0.2 * position_weight</code>. Apply position-aware sigmoid: <code>centered_ratio = ratio - 1.0</code>, <code>bounded_ratio = 1.0 + (2.0 * torch.sigmoid(centered_ratio) - 1.0) * effective_deviation</code>. Use: <code>surr1 = bounded_ratio * advantages</code>, <code>surr2 = torch.clamp(bounded_ratio, 1 - cliprange, 1 + cliprange) * advantages</code>, <code>loss = -torch.min(surr1, surr2)</code>. Add metadata: <code>metadata["mean_effective_deviation"] = effective_deviation.mean().item()</code>, <code>metadata["early_deviation"] = effective_deviation[:, :seq_len//4].mean().item()</code>, <code>metadata["late_deviation"] = effective_deviation[:, -seq_len//4:].mean().item()</code>.</p> <p><b>Validation Accuracy: 59.8</b></p>	<p><b>[Experiment]</b> Hierarchical Position-Group Trust Region: Apply trust region at two hierarchical levels – group level (shared within each response group) and position level (varies along sequence). Groups with high internal reward variance get tighter group-level bounds. Within groups, positions follow the proven decay pattern. This captures both cross-sample and within-sample structure. Formula: <math>\text{group\_dev} = 0.4 - 0.15 * \tanh(\text{group\_reward\_var} / 0.3)</math>, <math>\text{position\_factor} = 1 - 0.2 * \text{rel\_pos}</math>, <math>\text{effective\_dev} = \text{group\_dev} * \text{position\_factor}</math>.</p> <p><b>[Code Changes]</b> Modify <code>grpo.py</code>: Initialize <code>current_cliprange = 0.2</code>, <code>ema_clip_fraction = 0.15</code>. Standard momentum clip updates. Pass <code>group_size</code> to function. Modify <code>compute_grpo_clip_loss</code> in <code>grpo_utils.py</code>: Add parameter <code>group_size=8</code>. After computing ratio: <code>batch_size, seq_len = ratio.shape</code>, <code>n_groups = batch_size // group_size</code>. Compute group reward variance from advantages as proxy: <code>adv_grouped = advantages.view(n_groups, group_size, -1)</code>, <code>group_adv_var = adv_grouped.var(dim=1, keepdim=True)</code>, <code>group_adv_var_expanded = group_adv_var.expand(-1, group_size, -1).reshape(advantages.shape)</code>. Group-level deviation: <code>group_deviation = 0.4 - 0.15 * torch.tanh(group_adv_var_expanded / 0.3)</code>. Position factor: <code>positions = torch.arange(seq_len, device=ratio.device).float().unsqueeze(0).expand(batch_size, -1)</code>, <code>rel_pos = positions / (seq_len - 1 + 1e-6)</code>, <code>position_factor = 1.0 - 0.2 * rel_pos</code>. Hierarchical deviation: <code>effective_deviation = group_deviation * position_factor</code>, <code>effective_deviation = torch.clamp(effective_deviation, 0.15, 0.45)</code>. Apply: <code>centered_ratio = ratio - 1.0</code>, <code>bounded_ratio = 1.0 + (2.0 * torch.sigmoid(centered_ratio) - 1.0) * effective_deviation</code>. Use: <code>surr1 = bounded_ratio * advantages</code>, <code>surr2 = torch.clamp(bounded_ratio, 1 - cliprange, 1 + cliprange) * advantages</code>, <code>loss = -torch.min(surr1, surr2)</code>. Add metadata: <code>metadata["mean_group_var"] = group_adv_var.mean().item()</code>, <code>metadata["mean_effective_deviation"] = effective_deviation.mean().item()</code>. Log to wandb.</p>

Table 4: Additional examples on the GRPO environment. Ideas are generated by Claude-4.5-Opus during evolutionary search.

While the best-performing ideas on nanoGPT tend to be heavily optimized with extensive hyperparameter tuning mixed with various architecture tweaks, we pick a few more “atomic” algorithmic ideas that are successfully executed.

#### Examples from Claude-4.5-Opus

- 
- **Head-Wise Attention Output Scaling** Add learnable per-head scaling factors to attention, allowing different heads to contribute with different magnitudes to the output.  
**Validation Loss: 3.2386**
  - **Learned Residual Connection Weights** Add learnable scalar weights for each residual connection that are initialized to 1.0, allowing the model to learn optimal residual scaling during training.  
**Validation Loss: 3.2517**
  - **Mixture of Embeddings with Position** Learn to mix token embeddings and position embeddings with a content-dependent weight, allowing the model to dynamically balance positional vs semantic information per token.  
**Validation Loss: 3.2497**
  - **Shared Input-Output Embedding with Learned Asymmetry** Keep weight tying but add a small learned transformation on the output side, providing the benefits of weight tying while allowing output-specific adaptation.  
**Validation Loss: 3.2499**
  - **Gated Final Normalization** Replace the final RMSNorm before `lm_head` with a gated version where a learned gate controls how much normalization is applied vs passing the raw representation.  
**Validation Loss: 3.2503**
  - **Position-Aware MLP Gating** Gate the MLP output based on position information, allowing the model to learn position-dependent processing depth.  
**Validation Loss: 3.2506**
  - **Learned Residual Connection Weights** Add learnable scalar weights for each residual connection that are initialized to 1.0, allowing the model to learn optimal residual scaling during training.  
**Validation Loss: 3.2517**
  - **Grouped Token Embeddings** Group the vocabulary into clusters and add a learned embedding per cluster on top of token embeddings, providing hierarchical vocabulary structure.  
**Validation Loss: 3.2521**

Lastly, we present several interesting ideas on the nanoGPT environment that didn't get successfully executed. These examples are generated by Claude-4.5-Opus.

- **Soft Layer Repetition** Allow the model to softly repeat computation through layers by adding a learned gate that mixes the current layer's input back into its output, simulating variable depth.
- **Causal Context Compression** Before each attention layer, apply a learned compression that mixes local context (previous 2-3 tokens) into the current representation, providing implicit local context without convolutions.
- **Attention Head Specialization via Orthogonal Loss** Add a soft penalty that encourages different attention heads to attend to different patterns by penalizing similarity between head outputs.
- **Skip Connections with Learned Residual Weights** Combine skip connections with learned residual weights. The skip connections provide alternative gradient paths while learned weights allow adaptive scaling.
- **Token Difficulty-Aware Loss Weighting** Weight the loss contribution of each token based on the model's uncertainty (entropy) at that position, focusing learning on difficult tokens while not over-optimizing easy ones.

## A.5 CODE EXECUTION EXAMPLES

We present a few ideas with their full code execution to demonstrate the full end-to-end trajectories. All examples below are from Claude-4.5-Sonnet on the GRPO environment. For each example, we first present the natural language idea, followed by the code implementation generated by Claude-4.5-Sonnet.

### Example 1

**[Experiment]** Create mathematical working memory simulation by maintaining a context buffer of mathematical facts, definitions, and intermediate results during problem solving. This buffer gets updated as the model works through problems and provides additional context for subsequent mathematical steps, simulating how humans maintain mathematical working memory during complex calculations.

**[Code Changes]** Modify the prompt template in `sample.py` to include a “Mathematical Context” section that accumulates relevant mathematical facts during problem solving. Implement context buffer management in `train_loop` in `grpo.py` that maintains recent mathematical definitions, formulas, and intermediate calculations across problems in the same batch. Update the buffer with key mathematical expressions extracted from each response. Use enhanced prompts that include relevant context from the buffer. Set `-learning_rate 2.5e-5` and `-loss_type reinforce_with_baseline`.

**Validation Accuracy: 58.0**

#### Example 1 code diff

```
1 --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/
  sample.py
2 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/
  sample.py
3 @@ -2,6 +2,7 @@ from vllm import LLM, SamplingParams
4 import json
5 from tqdm import tqdm
6 import random
7 +import re
8
9 def load_prompt_template(prompt_path = "prompts/r1_zero.prompt"):
10
11     with open(prompt_path, 'r') as f:
12         @@ -20,12 +21,52 @@ def get_model_and_sampling_params(model_path
13             = "Qwen/Qwen2.5-Math-1.5B", token
14
15             return llm, sampling_params
16
17 -def load_dataset_with_prompt_template(prompt_template,
18     dataset_path="/juice5b/scr5b/nlp/aihinton/CS336_data/a5-
19     alignment/MATH/test.jsonl"):
20 +def extract_mathematical_expressions(text):
21 +    """Extract mathematical expressions, formulas, and key facts
22     from text."""
23 +    expressions = []
24 +
25 +    # Extract LaTeX math expressions
26 +    latex_patterns = [
27 +        r'\$[^$]+\$', # Inline math
28 +        r'\\\[a-zA-Z]+\{[^}]+\}', # LaTeX commands
29 +        r'\\frac\{[^}]+\}\{[^}]+\}', # Fractions
30 +        r'\\sqrt\{[^}]+\}', # Square roots
31 +        r'[a-zA-Z]\s*=\s*[^\n\.]+', # Variable assignments
32 +    ]
```

```

28 +
29 +     for pattern in latex_patterns:
30 +         matches = re.findall(pattern, text)
31 +         expressions.extend(matches)
32 +
33 +     # Extract definitions and key facts
34 +     definition_patterns = [
35 +         r'[A-Z][a-z]+ is defined as [^\.]+\.',
36 +         r'Let [a-zA-Z] be [^\.]+\.',
37 +         r'Given that [^\.]+\.',
38 +         r'Since [^\.]+\.',
39 +         r'Therefore [^\.]+\.',
40 +     ]
41 +
42 +     for pattern in definition_patterns:
43 +         matches = re.findall(pattern, text)
44 +         expressions.extend(matches)
45 +
46 +     return list(set(expressions)) # Remove duplicates
47 +
48 +def create_context_prompt(base_template, context_items):
49 +     """Create a prompt with mathematical context."""
50 +     if context_items:
51 +         context_section = "Mathematical Context:\n" + "\n".join(
52 +             f"- {item}" for item in context_items[-10:]) + "\n\n"
53 +         return context_section + base_template
54 +     return base_template
55 +def load_dataset_with_prompt_template(prompt_template,
56 +     dataset_path="/juice5b/scr5b/nlp/aihinton/CS336_data/a5-
57 +     alignment/MATH/test.jsonl", context_items=None):
58 +     with open(dataset_path, 'r') as f:
59 +         dataset = [json.loads(line) for line in f]
60 +         print(f"Loaded {len(dataset)} examples from {dataset_path}")
61 +     prompts = [prompt_template.format(question=example["problem
62 +     "]) for example in dataset]
63 +     enhanced_template = create_context_prompt(prompt_template,
64 +     context_items)
65 +     prompts = [enhanced_template.format(question=example["
66 +     problem"]) for example in dataset]
67 +     answers = [example["answer"] for example in dataset]
68 +     return prompts, answers
69 +
70 +@@ -93,4 +134,42 @@
71 +     if return_rewards:
72 +         return batched_prompts, batched_answers,
73 +         batched_responses, batched_rewards
74 +     return results
75 +
76 +
77 +class MathContextBuffer:
78 +     """Buffer to maintain mathematical working memory across
79 +     problems."""
80 +
81 +     def __init__(self, max_size=50):
82 +         self.buffer = []
83 +         self.max_size = max_size
84 +
85 +     def add_expressions(self, expressions):
86 +         """Add new mathematical expressions to the buffer."""
87 +         for expr in expressions:

```

```

82 +         if expr and len(expr.strip()) > 3: # Filter out
very short expressions
83 +             self.buffer.append(expr.strip())
84 +
85 +         # Keep only the most recent items
86 +         if len(self.buffer) > self.max_size:
87 +             self.buffer = self.buffer[-self.max_size:]
88 +
89 +     def get_relevant_context(self, query_text, max_items=10):
90 +         """Get the most relevant context items for a given query
. """
91 +         if not self.buffer:
92 +             return []
93 +
94 +         # Simple relevance scoring based on shared mathematical
terms
95 +         scored_items = []
96 +         query_lower = query_text.lower()
97 +
98 +         for item in self.buffer[-20:]: # Consider recent items
99 +             score = sum(1 for word in item.lower().split() if
word in query_lower)
100 +             if score > 0:
101 +                 scored_items.append((score, item))
102 +
103 +         # Return top items sorted by relevance
104 +         scored_items.sort(reverse=True, key=lambda x: x[0])
105 +         return [item[1] for item in scored_items[:max_items]]
106 +
107 +     def clear(self):
108 +         """Clear the context buffer."""
109 +         self.buffer = []
110 --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/grpo.
py
111 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/grpo.
py
112 @@ -1,6 +1,6 @@
113 from transformers import AutoModelForCausalLM, AutoTokenizer
114 from utils import tokenize_prompt_and_output,
get_response_log_probs
115 -from sample import load_prompt_template,
load_dataset_with_prompt_template, sample_rollout
116 +from sample import load_prompt_template,
load_dataset_with_prompt_template, sample_rollout,
MathContextBuffer, extract_mathematical_expressions,
create_context_prompt
117 from drgrpo_grader import rl_zero_reward_fn_train
118 from evaluate import rl_zero_reward_fn_eval, evaluate_vllm
119 from grpo_utils import compute_group_normalized_rewards,
grpo_microbatch_train_step
120 @@ -66,6 +66,9 @@ def evaluate_model(policy_model, vllm_model,
eval_prompts, eval_answers, eval_s
121 def train_loop(model, train_prompts, train_answers,
learning_rate, grpo_steps, train_steps_per_rollout,
output_dir, batch_size, gradient_accumulation_steps = 4,
group_size = 2, rollout_subset_size = 256, device = "cuda",
logging_steps = 20, saving_steps = 4000, eval_epochs = 5,
eval_prompts = None, eval_answers = None, sampling_params =
None, eval_vllm_model = None, cliprange = 0.2, loss_type = "
reinforce_with_baseline"):
122     model.to(device)
123     training_steps = grpo_steps
124 +

```

```

125 + # Initialize mathematical context buffer
126 + context_buffer = MathContextBuffer(max_size=100)
127 + optimizer = torch.optim.AdamW(model.parameters(), lr=
learning_rate, weight_decay=0.0, betas=(0.9, 0.95))
128 + global_step = 0 # Initialize global step counter
129
130 @@ -85,8 +88,31 @@ def train_loop(model, train_prompts,
train_answers, learning_rate, grpo_steps,
131 + load_policy_into_vllm_instance(model, vllm_model)
132
133 + ## sample rollouts
134 + # Get enhanced prompts with mathematical context for
this epoch
135 + print ("Sampling rollouts for epoch: ", epoch)
136 - rollout_prompts, rollout_answers, rollout_responses,
rollout_rewards = sample_rollout(vllm_model,
rl_zero_reward_fn_train, train_prompts, train_answers, G=
group_size, eval_sampling_params=eval_sampling_params,
subset_size=rollout_subset_size, return_rewards=True,
batch_size=512)
137 +
138 + # Create enhanced prompts with context for this batch
139 + enhanced_prompts = []
140 + base_template = load_prompt_template()
141 +
142 + # Select subset of problems for this epoch
143 + if rollout_subset_size is not None:
144 + indices = random.sample(range(len(train_prompts)),
rollout_subset_size)
145 + epoch_prompts = [train_prompts[i] for i in indices]
146 + epoch_answers = [train_answers[i] for i in indices]
147 + else:
148 + epoch_prompts = train_prompts
149 + epoch_answers = train_answers
150 +
151 + # Create context-enhanced prompts
152 + for prompt in epoch_prompts:
153 + relevant_context = context_buffer.
get_relevant_context(prompt, max_items=8)
154 + enhanced_prompt = create_context_prompt(
base_template.format(question=prompt.split("Question: ")[-1]),
relevant_context)
155 + enhanced_prompts.append(enhanced_prompt)
156 +
157 + # Sample with enhanced prompts (need to adapt
sample_rollout for direct prompt input)
158 + rollout_prompts, rollout_answers, rollout_responses,
rollout_rewards = sample_rollout_with_prompts(vllm_model,
rl_zero_reward_fn_train, enhanced_prompts, epoch_answers, G=
group_size, eval_sampling_params=eval_sampling_params,
return_rewards=True, batch_size=512)
159 +
160 + # Randomly sample 2 rollouts to print
161 + indices = random.sample(range(len(rollout_prompts)), 2)
162 + print ("Example rollouts:")
163 @@ -96,6 +122,15 @@ def train_loop(model, train_prompts,
train_answers, learning_rate, grpo_steps,
164 + print(f"Response: {rollout_responses[idx]}")
165 + print(f"Reward: {rollout_rewards[idx]}")
166 + print(f"Ground truth: {rollout_answers[idx]}")
167 +
168 + # Update context buffer with mathematical expressions
from responses

```

```

169 +     print("Updating mathematical context buffer...")
170 +     for response in rollout_responses:
171 +         if rollout_rewards[rollout_responses.index(response)
172 + ] > 0: # Only add from correct responses
173 +             mathematical_expressions =
174 + extract_mathematical_expressions(response)
175 +             if mathematical_expressions:
176 +                 context_buffer.add_expressions(
177 + mathematical_expressions)
178 +
179 + rollout_tokenized = tokenize_prompt_and_output(
180 + rollout_prompts, rollout_responses, tokenizer)
181 + rollout_data_loader = create_data_loader(
182 + rollout_tokenized, batch_size=batch_size, shuffle=False)
183 +
184 + @ -196,6 +231,47 @@ def train_loop(model, train_prompts,
185 + train_answers, learning_rate, grpo_steps,
186 + ), step=global_step)
187 +     # print ("Global Step: ", global_step, "Loss:
188 + ", loss.item(), "Entropy: ", entropy.mean().item(), "Clip
189 + fraction: ", metadata.get("clip_fraction", 0.0))
190 +
191 + def sample_rollout_with_prompts(
192 + vllm_model,
193 + reward_fn,
194 + prompts,
195 + answers,
196 + G,
197 + eval_sampling_params,
198 + return_rewards=False,
199 + batch_size=64
200 + ):
201 +     """Sample rollouts using provided prompts directly."""
202 +     # Create batched prompts by repeating each prompt G times
203 +     batched_prompts = []
204 +     batched_answers = []
205 +     for prompt, answer in zip(prompts, answers):
206 +         batched_prompts.extend([prompt] * G)
207 +         batched_answers.extend([answer] * G)
208 +
209 +     # Process in batches to avoid OOM
210 +     all_outputs = []
211 +     for i in range(0, len(batched_prompts), batch_size):
212 +         batch_prompts = batched_prompts[i:i + batch_size]
213 +         batch_outputs = vllm_model.generate(batch_prompts,
214 + eval_sampling_params)
215 +         all_outputs.extend(batch_outputs)
216 +
217 +     # Process results
218 +     batched_responses = []
219 +     batched_rewards = []
220 +     total_rewards = 0
221 +
222 +     for output, answer in tqdm(zip(all_outputs, batched_answers)
223 + ):
224 +         generated_answer = output.outputs[0].text
225 +         reward = reward_fn(generated_answer, answer)
226 +         batched_responses.append(generated_answer)
227 +         batched_rewards.append(reward["reward"])
228 +         total_rewards += reward["reward"]
229 +

```

```

220 +     print(f"Accuracy of sampled rollouts: {total_rewards}/{len(
      batched_prompts)} = {total_rewards / len(batched_prompts) *
      100}%")
221 +
222 +     return batched_prompts, batched_answers, batched_responses,
      batched_rewards
223 +
224 + if __name__ == "__main__":
225     parser = argparse.ArgumentParser()
226     parser.add_argument("--model_path", type=str, default="Qwen/
      Qwen2.5-Math-1.5B")
227 @@ -203,7 +279,7 @@ if __name__ == "__main__":
228     parser.add_argument("--train_dataset_path", type=str,
      default="../MATH/train.jsonl")
229     parser.add_argument("--eval_dataset_path", type=str, default
      = "../MATH/test.jsonl")
230     parser.add_argument("--output_dir", type=str, default="ckpts
      /")
231 -     parser.add_argument("--learning_rate", type=float, default=1
      e-5)
232 +     parser.add_argument("--learning_rate", type=float, default
      =2.5e-5)
233     parser.add_argument("--grpo_steps", type=int, default=200)
234     parser.add_argument("--group_size", type=int, default=8)
235     parser.add_argument("--rollout_subset_size", type=int,
      default=256)
236 @@ -212,7 +288,7 @@ if __name__ == "__main__":
237     parser.add_argument("--gradient_accumulation_steps", type=
      int, default=16)
238     parser.add_argument("--batch_size", type=int, default=4)
239     parser.add_argument("--cliprange", type=float, default=0.2)
240 -     parser.add_argument("--loss_type", type=str, default="
      grpo_clip")
241 +     parser.add_argument("--loss_type", type=str, default="
      reinforce_with_baseline")
242     parser.add_argument("--wandb_project", type=str, default="
      grpo-math")
243     parser.add_argument("--wandb_name", type=str, default="
      grpo_clip_1")
244     args = parser.parse_args()
245 --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/
      run_job.sh
246 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_43/
      run_job.sh
247 @@ -21,7 +21,7 @@ timeout 2h uv run \
248     --index https://download.pytorch.org/whl/cul28 \
249     --index-strategy unsafe-best-match \
250     python grpo.py \
251 -     --learning_rate 1e-5 \
252 +     --learning_rate 2.5e-5 \
253     --grpo_steps 20 \
254     --group_size 8 \
255     --rollout_subset_size 128 \
256 @@ -30,7 +30,7 @@ timeout 2h uv run \
257     --gradient_accumulation_steps 16 \
258     --batch_size 4 \
259     --cliprange 0.2 \
260 -     --loss_type grpo_clip \
261 +     --loss_type reinforce_with_baseline \
262     --wandb_name $wandb_name
263
264     echo "Experiment finished successfully!"

```

## Example 2

**[Experiment]** Implement mathematical solution robustness training through systematic perturbation testing that teaches the model to solve mathematically equivalent problems with varied presentations, notation styles, and problem phrasings. Combine this with proven reward shaping by creating robustness-aware rewards that encourage mathematical understanding that generalizes across problem variations.

**[Code Changes]** Create problem perturbation system in `sample.py` that generates equivalent mathematical problems with varied notation, different variable names, alternative problem phrasing, and equivalent mathematical formulations. Track solution consistency across perturbations. Enhance `rl_zero_reward_fn_train` in `drgrp0_grader.py` to reward robustness: `robustness_bonus = 0.12 * consistency_score_across_perturbations` where consistency measures solution method stability across equivalent problems. Apply generalization bonuses: `+0.08` for solving problems with unfamiliar notation, `+0.06` for handling varied problem presentations. Modify `train_loop` in `grp0.py` to include perturbation-based training batches and track robustness improvement. Use robustness-aware learning rate with cosine annealing from `3.4e-5` to `1.8e-5`. Implement perturbation-complexity gradient clipping. Set `-loss_type reinforce_with_baseline`.

## Validation Accuracy: 56.2

### Example 2 code diff

```
1 --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/
    sample.py
2 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/
    sample.py
3 @@ -2,6 +2,7 @@ from vllm import LLM, SamplingParams
4 import json
5 from tqdm import tqdm
6 import random
7 +import re
8
9 def load_prompt_template(prompt_path = "prompts/rl_zero.prompt"):
10     with open(prompt_path, 'r') as f:
11 @@ -29,6 +30,125 @@ def load_dataset_with_prompt_template(
    prompt_template, dataset_path="/juice5b/s
    answers = [example["answer"] for example in dataset]
13     return prompts, answers
14
15 +def create_notation_perturbations(problem_text):
16 +     """Create mathematical notation variations"""
17 +     perturbations = []
18 +
19 +     # Variable substitutions
20 +     var_substitutions = [
21 +         ('x', 'a'), ('y', 'b'), ('z', 'c'),
22 +         ('n', 'm'), ('k', 'j'), ('p', 'q'),
23 +         ('r', 's'), ('t', 'u'), ('v', 'w')
24 +     ]
25 +
26 +     # Create variable name variations
27 +     for old_var, new_var in var_substitutions:
28 +         if old_var in problem_text:
29 +             new_text = problem_text.replace(old_var, new_var)
30 +             if new_text != problem_text:
31 +                 perturbations.append(("notation_var", new_text))
32 +
```

```

33 + # Mathematical symbol variations
34 + symbol_substitutions = [
35 +     ('\\cdot', '\\times'),
36 +     ('\\times', '\\cdot'),
37 +     ('\\frac', '\\div'),
38 +     ('\\le', '\\leq'),
39 +     ('\\ge', '\\geq'),
40 + ]
41 +
42 + for old_sym, new_sym in symbol_substitutions:
43 +     if old_sym in problem_text:
44 +         new_text = problem_text.replace(old_sym, new_sym)
45 +         perturbations.append(("notation_symbol", new_text))
46 +
47 + return perturbations
48 +
49 +def create_phrasing_perturbations(problem_text):
50 +     """Create alternative problem phrasings"""
51 +     perturbations = []
52 +
53 +     # Common phrasing substitutions
54 +     phrasing_patterns = [
55 +         (r"Find the value of", "Determine"),
56 +         (r"What is", "Calculate"),
57 +         (r"Solve for", "Find"),
58 +         (r"How many", "What is the number of"),
59 +         (r"Compute", "Find"),
60 +         (r"Evaluate", "Calculate"),
61 +     ]
62 +
63 +     for pattern, replacement in phrasing_patterns:
64 +         if re.search(pattern, problem_text, re.IGNORECASE):
65 +             new_text = re.sub(pattern, replacement, problem_text,
66 +                               flags=re.IGNORECASE)
67 +             if new_text != problem_text:
68 +                 perturbations.append(("phrasing", new_text))
69 +
70 +     return perturbations
71 +def create_formulation_perturbations(problem_text):
72 +     """Create equivalent mathematical formulations"""
73 +     perturbations = []
74 +
75 +     # Simple algebraic reformulations
76 +     reformulations = [
77 +         (r"(\w+) = (\d+) \+ (\d+)", r"\1 - \2 = \3"), # a = b +
78 +         (r"(\w+) \+ (\d+) = (\d+)", r"\1 = \3 - \2"), # x + a =
79 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
80 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
81 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
82 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
83 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
84 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
85 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
86 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
87 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
88 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =
89 +         (r"(\d+) - (\w+) = (\d+)", r"\2 = \1 - \3"), # a - x =

```

```

90 +def generate_problem_perturbations(problem_text,
91 +     max_perturbations=3):
92 +     """Generate various perturbations of a mathematical problem
93 +     """
94 +     all_perturbations = []
95 +     # Generate different types of perturbations
96 +     all_perturbations.extend(create_notation_perturbations(
97 +         problem_text))
98 +     all_perturbations.extend(create_phrasing_perturbations(
99 +         problem_text))
100 +     all_perturbations.extend(create_formulation_perturbations(
101 +         problem_text))
102 +     # Randomly sample up to max_perturbations
103 +     if len(all_perturbations) > max_perturbations:
104 +         all_perturbations = random.sample(all_perturbations,
105 +             max_perturbations)
106 +     return all_perturbations
107 +
108 +def sample_with_perturbations(prompts, answers,
109 +     perturbation_ratio=0.3):
110 +     """Sample problems with perturbations mixed in"""
111 +     perturbed_prompts = []
112 +     perturbed_answers = []
113 +     perturbation_metadata = []
114 +
115 +     for prompt, answer in zip(prompts, answers):
116 +         # Always include original
117 +         perturbed_prompts.append(prompt)
118 +         perturbed_answers.append(answer)
119 +         perturbation_metadata.append({"type": "original", "
120 +             source_idx": len(perturbed_prompts)-1})
121 +
122 +         # Add perturbations with some probability
123 +         if random.random() < perturbation_ratio:
124 +             # Extract problem text from prompt (assuming it's in
125 +             a specific format)
126 +             problem_text = prompt.split("question:")[-1] if "
127 +             question:" in prompt else prompt
128 +             perturbations = generate_problem_perturbations(
129 +                 problem_text)
130 +
131 +             for pert_type, pert_text in perturbations:
132 +                 # Reconstruct full prompt with perturbed problem
133 +                 pert_prompt = prompt.replace(problem_text,
134 +                     pert_text)
135 +                 perturbed_prompts.append(pert_prompt)
136 +                 perturbed_answers.append(answer) # Same answer
137 +                 for equivalent problem
138 +                     perturbation_metadata.append({
139 +                         "type": pert_type,
140 +                         "source_idx": len(perturbed_prompts)-len(
141 +                             perturbations)-1
142 +                     })
143 +
144 +     return perturbed_prompts, perturbed_answers,
145 +     perturbation_metadata
146 +
147 +def sample_rollout(
148 +     vllm_model,

```

```

137 --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/
138 drgrpo_grader.py
139 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/
140 drgrpo_grader.py
141 @@ -1025,3 +1025,83 @@ def rl_zero_reward_fn_train(response,
142 ground_truth, fast=True):
143     "reward": 0.0
144     }
145
146 +def compute_consistency_score(responses, ground_truth,
147 perturbation_metadata):
148 +    """Compute consistency score across perturbations"""
149 +    if not perturbation_metadata:
150 +        return 0.0
151 +
152 +    # Group responses by source problem
153 +    source_groups = {}
154 +    for i, meta in enumerate(perturbation_metadata):
155 +        source_idx = meta.get("source_idx", i)
156 +        if source_idx not in source_groups:
157 +            source_groups[source_idx] = []
158 +        source_groups[source_idx].append((i, responses[i], meta))
159 +
160 +    total_consistency = 0.0
161 +    valid_groups = 0
162 +
163 +    for source_idx, group_data in source_groups.items():
164 +        if len(group_data) > 1: # Only consider groups with
165 +            multiple responses
166 +                # Get correctness for each response in the group
167 +                correctness_scores = []
168 +                for _, response, meta in group_data:
169 +                    reward_result = rl_zero_reward_fn_train(response,
170 ground_truth)
171 +                    correctness_scores.append(reward_result["reward
172 "])
173 +
174 +                # Consistency is measured as how often the model
175 +                gets the same correctness
176 +                if correctness_scores:
177 +                    # If all responses have same correctness (all
178 +                    correct or all incorrect), high consistency
179 +                    unique_scores = set(correctness_scores)
180 +                    if len(unique_scores) == 1:
181 +                        consistency = 1.0
182 +                    else:
183 +                        # Partial consistency based on agreement
184 +                        most_common = max(set(correctness_scores),
185 key=correctness_scores.count)
186 +                        consistency = correctness_scores.count(
187 most_common) / len(correctness_scores)
188 +
189 +                        total_consistency += consistency
190 +                        valid_groups += 1
191 +
192 +    if valid_groups == 0:
193 +        return 0.0
194 +
195 +    return total_consistency / valid_groups
196 +
197 +def rl_zero_reward_fn_train_robust(response, ground_truth,
198 perturbation_metadata=None, responses_batch=None, fast=True):

```

```

187 +     """Enhanced reward function with robustness bonuses"""
188 +     # Get base reward
189 +     base_reward = rl_zero_reward_fn_train(response, ground_truth,
190 +     fast)
191 +
192 +     # Initialize bonus components
193 +     robustness_bonus = 0.0
194 +     notation_bonus = 0.0
195 +     presentation_bonus = 0.0
196 +
197 +     # Compute robustness bonus if perturbation data available
198 +     if perturbation_metadata and responses_batch:
199 +         consistency_score = compute_consistency_score(
200 +         responses_batch, ground_truth, perturbation_metadata)
201 +         robustness_bonus = 0.12 * consistency_score
202 +
203 +     # Check for notation and presentation bonuses from metadata
204 +     if perturbation_metadata:
205 +         current_meta = perturbation_metadata[0] if
206 +         perturbation_metadata else {}
207 +         pert_type = current_meta.get("type", "original")
208 +
209 +         if pert_type in ["notation_var", "notation_symbol"] and
210 +         base_reward["reward"] > 0:
211 +             notation_bonus = 0.08
212 +
213 +         if pert_type in ["phrasing", "formulation"] and
214 +         base_reward["reward"] > 0:
215 +             presentation_bonus = 0.06
216 +
217 +     # Combine all rewards
218 +     total_reward = base_reward["reward"] + robustness_bonus +
219 +     notation_bonus + presentation_bonus
220 +
221 +     return {
222 +         "format_reward": base_reward["format_reward"],
223 +         "answer_reward": base_reward["answer_reward"],
224 +         "robustness_bonus": robustness_bonus,
225 +         "notation_bonus": notation_bonus,
226 +         "presentation_bonus": presentation_bonus,
227 +         "reward": total_reward
228 +     }
229 +
230 + --- repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/grpo.
231 + py
232 + +++ repo_variants_full_pipeline_bsz50_fixed2_epoch7/idea_37/grpo.
233 + py
234 + @@ -1,7 +1,7 @@
235 + from transformers import AutoModelForCausalLM, AutoTokenizer
236 + from utils import tokenize_prompt_and_output,
237 +     get_response_log_probs
238 + -from sample import load_prompt_template,
239 +     load_dataset_with_prompt_template, sample_rollout
240 + -from drgrpo_grader import rl_zero_reward_fn_train
241 + +from sample import load_prompt_template,
242 +     load_dataset_with_prompt_template, sample_rollout,
243 +     sample_with_perturbations
244 + +from drgrpo_grader import rl_zero_reward_fn_train,
245 +     rl_zero_reward_fn_train_robust
246 + from evaluate import rl_zero_reward_fn_eval, evaluate_vllm
247 + from grpo_utils import compute_group_normalized_rewards,
248 +     grpo_microbatch_train_step
249 + from torch.utils.data import DataLoader, Dataset
250 + @@ -12,6 +12,7 @@ from tqdm import tqdm

```

```

236 from vllm import LLM, SamplingParams
237 import wandb
238 import random
239 +import math
240
241 def load_policy_into_vllm_instance(policy, llm):
242     state_dict = policy.state_dict()
243 @@ -63,11 +64,23 @@
244     metrics = evaluate_vllm(vllm_model, rl_zero_reward_fn_eval,
245                             eval_prompts, eval_answers, eval_sampling_params, output_path
246                             =output_path)
247     return metrics
248
249 -def train_loop(model, train_prompts, train_answers,
250                learning_rate, grpo_steps, train_steps_per_rollout,
251                output_dir, batch_size, gradient_accumulation_steps = 4,
252                group_size = 2, rollout_subset_size = 256, device = "cuda",
253                logging_steps = 20, saving_steps = 4000, eval_epochs = 5,
254                eval_prompts = None, eval_answers = None, sampling_params =
255                None, eval_vllm_model = None, cliprange = 0.2, loss_type = "
256                reinforce_with_baseline"):
257 +def get_cosine_annealing_lr(step, total_steps, lr_max=3.4e-5,
258                             lr_min=1.8e-5):
259 +     """Cosine annealing learning rate schedule"""
260 +     return lr_min + (lr_max - lr_min) * 0.5 * (1 + math.cos(math.
261 +         pi * step / total_steps))
262 +
263 +def compute_perturbation_complexity_clip_norm(model,
264         complexity_factor=1.0):
265 +     """Compute gradient clipping norm based on perturbation
266 +     complexity"""
267 +     base_clip_norm = 1.0
268 +     return base_clip_norm * complexity_factor
269 +
270 +def train_loop(model, train_prompts, train_answers,
271                learning_rate, grpo_steps, train_steps_per_rollout,
272                output_dir, batch_size, gradient_accumulation_steps = 4,
273                group_size = 2, rollout_subset_size = 256, device = "cuda",
274                logging_steps = 20, saving_steps = 4000, eval_epochs = 5,
275                eval_prompts = None, eval_answers = None, sampling_params =
276                None, eval_vllm_model = None, cliprange = 0.2, loss_type = "
277                reinforce_with_baseline", perturbation_ratio = 0.3):
278     model.to(device)
279     training_steps = grpo_steps
280     optimizer = torch.optim.AdamW(model.parameters(), lr=
281     learning_rate, weight_decay=0.0, betas=(0.9, 0.95))
282     # Start with robustness-aware learning rate
283     initial_lr = get_cosine_annealing_lr(0, grpo_steps)
284     optimizer = torch.optim.AdamW(model.parameters(), lr=
285     initial_lr, weight_decay=0.0, betas=(0.9, 0.95))
286     global_step = 0 # Initialize global step counter
287     robustness_scores = []
288
289     for epoch in range(grpo_steps):
290         model.train()
291 @@ -82,21 +95,50 @@ def train_loop(model, train_prompts,
292     train_answers, learning_rate, grpo_steps,
293
294     model.train()
295
296     # Update learning rate with cosine annealing
297     current_lr = get_cosine_annealing_lr(epoch, grpo_steps)
298     for param_group in optimizer.param_groups:

```

```

276 +         param_group['lr'] = current_lr
277 +
278 +         ## load the current policy model to vllm for sampling
rollouts
279 +         load_policy_into_vllm_instance(model, vllm_model)
280 +
281 +         # Sample with perturbations for robustness training
282 +         perturbed_prompts, perturbed_answers,
perturbation_metadata = sample_with_perturbations(
283 +             train_prompts, train_answers, perturbation_ratio=
perturbation_ratio
284 +         )
285 +
286 +         ## sample rollouts
287 +         print ("Sampling rollouts for epoch: ", epoch)
288 -         rollout_prompts, rollout_answers, rollout_responses,
rollout_rewards = sample_rollout(vllm_model,
rl_zero_reward_fn_train, train_prompts, train_answers, G=
group_size, eval_sampling_params=eval_sampling_params,
subset_size=rollout_subset_size, return_rewards=True,
batch_size=512)
289 +
290 +         # Use subset of perturbed prompts for training
291 +         subset_size = min(rollout_subset_size, len(
perturbed_prompts))
292 +         if subset_size < len(perturbed_prompts):
293 +             indices = random.sample(range(len(perturbed_prompts)
), subset_size)
294 +             subset_prompts = [perturbed_prompts[i] for i in
indices]
295 +             subset_answers = [perturbed_answers[i] for i in
indices]
296 +             subset_metadata = [perturbation_metadata[i] for i in
indices]
297 +         else:
298 +             subset_prompts = perturbed_prompts
299 +             subset_answers = perturbed_answers
300 +             subset_metadata = perturbation_metadata
301 +
302 +             rollout_prompts, rollout_answers, rollout_responses,
rollout_rewards = sample_rollout(
303 +                 vllm_model, rl_zero_reward_fn_train, subset_prompts,
subset_answers,
304 +                 G=group_size, eval_sampling_params=
eval_sampling_params,
305 +                 subset_size=None, return_rewards=True, batch_size
=512
306 +             )
307 +
308 +         # Randomly sample 2 rollouts to print
309 +         indices = random.sample(range(len(rollout_prompts)), 2)
310 +         print ("Example rollouts:")
311 +         for idx in indices:
312 +             print(f"\nRollout {idx}:")
313 -             print(f"Prompt: {rollout_prompts[idx]}")
314 +             print(f"Prompt: {rollout_prompts[idx][:200]}...")
315 +             print(f"Response: {rollout_responses[idx]}")
316 +             print(f"Reward: {rollout_rewards[idx]}")
317 -             print(f"Ground truth: {rollout_answers[idx]}")
318 +             print(f"Ground truth: {rollout_answers[idx
][:100]}...")
319 +

```

```

320     rollout_tokenized = tokenize_prompt_and_output(
321         rollout_prompts, rollout_responses, tokenizer)
322     rollout_data_loader = create_data_loader(
323         rollout_tokenized, batch_size=batch_size, shuffle=False)
324
325     @@ -126,15 +168,32 @@ def train_loop(model, train_prompts,
326         train_answers, learning_rate, grpo_steps,
327
328         # Compute advantages using group normalization - no
329         # gradients needed
330         with torch.no_grad():
331             advantages, raw_rewards, metadata =
332             compute_group_normalized_rewards(
333                 reward_fn=rl_zero_reward_fn_train,
334                 # Create enhanced reward function with robustness
335                 def robust_reward_fn(response, ground_truth):
336                     # Find corresponding metadata for this response
337                     response_idx = rollout_responses.index(response)
338                     if response in rollout_responses else 0
339                     meta_idx = min(response_idx, len(subset_metadata
340                                 ) - 1)
341                     current_meta = [subset_metadata[meta_idx]] if
342                     subset_metadata else None
343
344                     return rl_zero_reward_fn_train_robust(
345                         response, ground_truth,
346                         perturbation_metadata=current_meta,
347                         responses_batch=rollout_responses
348                     )
349             advantages, raw_rewards, metadata =
350             compute_group_normalized_rewards(
351                 reward_fn=robust_reward_fn,
352                 rollout_responses=rollout_responses,
353                 repeated_ground_truths=rollout_answers,
354                 group_size=group_size,
355                 advantage_eps=1e-6,
356                 normalize_by_std=True
357             )
358             advantages = advantages.to(device)
359
360             # Track robustness improvement
361             current_robustness = metadata.get('mean_reward',
362             0.0)
363             robustness_scores.append(current_robustness)
364
365             # Log raw rewards statistics
366             print("\nGRPO epoch: ", epoch)
367             @@ -145,11 +204,20 @@
368             wandb.log({
369                 "eval/mean_reward": eval_mean_reward,
370                 "train/mean_reward": metadata["mean_reward"],
371                 "train/learning_rate": current_lr,
372                 "train/robustness_score": current_robustness,
373             }, step=global_step)
374             else:
375                 wandb.log({
376                     "train/mean_reward": metadata["mean_reward"],
377                     "train/learning_rate": current_lr,
378                     "train/robustness_score": current_robustness,
379                 }, step=global_step)
380
381             # Compute perturbation complexity for gradient clipping

```

```

373 +         perturbation_types = set(meta.get("type", "original")
374 +             for meta in subset_metadata)
375 +         complexity_factor = 1.0 + 0.1 * len(perturbation_types)
376 +         # More complex with more perturbation types
377 +         clip_norm = compute_perturbation_complexity_clip_norm(
378 +             model, complexity_factor)
379 @@ -185,6 +252,9 @@ def train_loop(model, train_prompts,
380 +             train_answers, learning_rate, grpo_steps,
381 +             )
382 +             if (batch_idx + 1) % gradient_accumulation_steps
383 +                 == 0:
384 +                     # Apply perturbation-complexity gradient
385 +                     clipping
386 +                     torch.nn.utils.clip_grad_norm_(model.
387 +                         parameters(), clip_norm)
388 +                     optimizer.step()
389 +                     optimizer.zero_grad()
390 @@ -211,7 +281,7 @@ if __name__ == "__main__":
391 +     parser.add_argument("--gradient_accumulation_steps", type=
392 +         int, default=16)
393 +     parser.add_argument("--batch_size", type=int, default=4)
394 +     parser.add_argument("--cliprange", type=float, default=0.2)
395 -     parser.add_argument("--loss_type", type=str, default="
396 +         grpo_clip")
397 +     parser.add_argument("--loss_type", type=str, default="
398 +         reinforce_with_baseline")
399 +     parser.add_argument("--wandb_project", type=str, default="
400 +         grpo-math")
401 +     parser.add_argument("--wandb_name", type=str, default="
402 +         grpo_clip_1")
403 +     args = parser.parse_args()
404 @@ -266,7 +336,8 @@ if __name__ == "__main__":
405 +     sampling_params=eval_sampling_params,
406 +     eval_vllm_model=vllm_model,
407 +     cliprange=args.cliprange,
408 +     loss_type=args.loss_type
409 -     loss_type=args.loss_type,
410 +     perturbation_ratio=0.3
411 +     )

```

### Example 3

**[Experiment]** Implement token-level reward attribution by using attention weights to identify which input tokens contributed most to correct answers, then amplifying the gradient updates for those tokens during policy gradient training.

**[Code Changes]** Modify `get_response_log_probs` in `utils.py` to also return attention weights from the last layer. In `grpo_microbatch_train_step`, compute token importance scores by averaging attention weights across heads, then multiply the policy gradient loss by  $(1 + \text{importance\_score})$  for tokens with high attention to mathematical content.

**Validation Accuracy: 45.2**

#### Example 3 code diff

```
1 --- repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/utils.
  py
2 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/utils.
  py
3 @@ -41,12 +41,13 @@ def compute_entropy(logits):
4     entropy = -torch.sum(probs * log_probs, dim=-1)
5     return entropy
6
7 -def get_response_log_probs(model, input_ids, labels,
  return_token_entropy=False, no_grad=True):
8 +def get_response_log_probs(model, input_ids, labels,
  return_token_entropy=False, return_attention=False, no_grad=
  True):
9     if no_grad:
10         with torch.no_grad():
11 -             outputs = model(input_ids, labels=labels)
12 +             outputs = model(input_ids, labels=labels,
  output_attentions=return_attention)
13             logits = outputs.logits # (batch_size, seq_len,
  vocab_size)
14             log_probs = torch.log_softmax(logits, dim=-1) # (
  batch_size, seq_len, vocab_size)
15 +             attentions = outputs.attentions if return_attention
  else None
16             # Get log probs of the actual label tokens
17             batch_size, seq_len = labels.shape # (batch_size,
  seq_len)
18             log_probs = torch.gather(log_probs, dim=-1, index=
  labels.unsqueeze(-1).squeeze(-1))
19 @@ -55,8 +56,9 @@ def get_response_log_probs(model, input_ids,
  labels, return_token_entropy=False
20             else:
21                 entropy = None
22     else:
23 -         outputs = model(input_ids, labels=labels)
24 +         outputs = model(input_ids, labels=labels,
  output_attentions=return_attention)
25             logits = outputs.logits # (batch_size, seq_len,
  vocab_size)
26 +             attentions = outputs.attentions if return_attention else
  None
27             log_probs = torch.log_softmax(logits, dim=-1) # (
  batch_size, seq_len, vocab_size)
28             # Get log probs of the actual label tokens
29             batch_size, seq_len = labels.shape # (batch_size,
  seq_len)
```

```

30 @@ -65,9 +67,17 @@ def get_response_log_probs(model, input_ids,
    labels, return_token_entropy=False
31         entropy = compute_entropy(logits)
32     else:
33         entropy = None
34 -
35 -     return {
36 +
37 +     result = {
38         "log_probs": log_probs,
39         "token_entropy": entropy
40     }
41 +
42 +     if return_attention and attentions is not None:
43 +         # Return attention weights from the last layer, averaged
44 +         # across heads
45 +         last_layer_attention = attentions[-1] # Shape: (
46 +         batch_size, num_heads, seq_len, seq_len)
47 +         averaged_attention = last_layer_attention.mean(dim=1) #
48 +         Average across heads: (batch_size, seq_len, seq_len)
49 +         result["attention_weights"] = averaged_attention
50 --- repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/
51 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/
52 @@ -170,6 +170,7 @@ def grpo_microbatch_train_step(
53     advantages: torch.Tensor | None = None,
54     old_log_probs: torch.Tensor | None = None,
55     cliprange: float | None = None,
56 +     attention_weights: torch.Tensor | None = None,
57 ) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
58     '''
59     Return:
60 @@ -180,6 +181,20 @@ def grpo_microbatch_train_step(
61     You should call loss.backward() in this function. Make sure
62     to adjust for gradient accumulation.
63     '''
64     loss, metadata = compute_policy_gradient_loss(
65     policy_log_probs, loss_type, raw_rewards, advantages,
66     old_log_probs, cliprange) # (batch_size, sequence_length)
67 +
68 +     # Apply attention-based weighting if attention weights are
69 +     # provided
70 +     if attention_weights is not None:
71 +         # Compute importance scores by averaging attention
72 +         # weights to earlier tokens
73 +         # attention_weights shape: (batch_size, seq_len, seq_len)
74 +         )
75 +         # For each token, sum attention weights to all previous
76 +         # tokens (including itself)
77 +         importance_scores = torch.sum(attention_weights, dim=-1)
78 +         # (batch_size, seq_len)
79 +
80 +         # Normalize importance scores to [0, 1] range per
81 +         # sequence
82 +         importance_scores = importance_scores / (
83 +         importance_scores.max(dim=-1, keepdim=True)[0] + 1e-8)
84 +
85 +         # Amplify loss for tokens with high importance: multiply
86 +         # by (1 + importance_score)

```

```

76 +         loss = loss * (1.0 + importance_scores)
77 +
78     loss = masked_mean(loss, response_mask)
79     loss = loss / gradient_accumulation_steps
80     loss.backward()
81 --- repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/grpo.
    py
82 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch1/idea_32/grpo.
    py
83 @@ -109,6 +109,7 @@ def train_loop(model, train_prompts,
    train_answers, learning_rate, grpo_steps,
84         model,
85         input_ids,
86         labels,
87         return_token_entropy=False,
88 +         return_attention=False,
89         no_grad=True
90     )
91 @@ -163,11 +164,13 @@
92         model,
93         input_ids,
94         labels,
95         return_token_entropy=True,
96 +         return_attention=True,
97         no_grad=False
98     )
99     policy_log_probs = response_log_probs["log_probs
100     "]
101     entropy = response_log_probs["token_entropy"]
102     attention_weights = response_log_probs.get("
103     attention_weights")
104
105     # Calculate data index for advantages/
106     old_log_probs
107     batch_idx_total = batch_idx * batch_size
108 @@ -177,7 +180,8 @@ def train_loop(model, train_prompts,
109     train_answers, learning_rate, grpo_steps,
110     loss_type=loss_type,
111     advantages=batch_advantages,
112     old_log_probs=batch_old_log_probs,
113     cliprange=cliprange
114     cliprange=cliprange,
115     attention_weights=attention_weights
116 )
117
118 if (batch_idx + 1) % gradient_accumulation_steps
119 == 0:

```

#### Example 4

**[Experiment]** Implement response diversity rewards within groups where responses to the same prompt receive bonus rewards (0.05-0.15) for being dissimilar to other responses in their group, encouraging exploration of different solution paths while maintaining the proven `group_size=8` and `3e-5` learning rate combination.

**[Code Changes]** Modify `compute_group_normalized_rewards` in `grp_utils.py` to compute pairwise similarity between responses in each group using token-level Jaccard similarity. Add diversity bonus: `diversity_reward = 0.15 * (1 - max_similarity_in_group)` to each response's reward before advantage computation. Reshape responses into groups, compute similarities, and add bonuses before advantage normalization. Set `-learning_rate 3e-5`, `-loss_type reinforce_with_baseline`, `-group_size 8`.

**Validation Accuracy: 19.2**

#### Example 4 code diff

```
1 --- repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/
    grp_utils.py
2 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/
    grp_utils.py
3 @@ -1,5 +1,25 @@
4   import torch
5   from typing import Literal
6
7 +def compute_jaccard_similarity(response1: str, response2: str)
   -> float:
8 +   """
9 +   Compute token-level Jaccard similarity between two responses.
10 +
11 +   Jaccard similarity = |intersection| / |union|
12 +   """
13 +   # Simple tokenization by splitting on whitespace
14 +   tokens1 = set(response1.split())
15 +   tokens2 = set(response2.split())
16 +
17 +   if len(tokens1) == 0 and len(tokens2) == 0:
18 +       return 1.0
19 +
20 +   intersection = tokens1.intersection(tokens2)
21 +   union = tokens1.union(tokens2)
22 +
23 +   if len(union) == 0:
24 +       return 1.0
25 +
26 +   return len(intersection) / len(union)
27
28 def compute_group_normalized_rewards(
29     reward_fn,
30     @@ -30,31 +49,62 @@
31     raw_rewards.append(reward_dict["reward"])
32     raw_rewards = torch.tensor(raw_rewards)
33
34 +   # Compute diversity bonuses within each group
35 +   n_groups = len(rollout_responses) // group_size
36 +   diversity_bonuses = []
37 +
38 +   for group_idx in range(n_groups):
39 +       group_start = group_idx * group_size
```

```

39 +         group_end = group_start + group_size
40 +         group_responses = rollout_responses[group_start:
group_end]
41 +
42 +         # Compute pairwise similarities within the group
43 +         for i in range(group_size):
44 +             max_similarity = 0.0
45 +             for j in range(group_size):
46 +                 if i != j:
47 +                     similarity = compute_jaccard_similarity(
group_responses[i], group_responses[j])
48 +                     max_similarity = max(max_similarity,
similarity)
49 +
50 +             # Diversity bonus: higher reward for more dissimilar
responses
51 +             diversity_bonus = 0.15 * (1 - max_similarity)
52 +             diversity_bonuses.append(diversity_bonus)
53 +
54 +         diversity_bonuses = torch.tensor(diversity_bonuses)
55 +
56 +         # Reshape rewards into groups
57 +         n_groups = len(raw_rewards) // group_size
58 +         grouped_rewards = raw_rewards.view(n_groups, group_size)
59 +         grouped_diversity_bonuses = diversity_bonuses.view(n_groups,
group_size)
60 +
61 +         # Add diversity bonuses to raw rewards before advantage
computation
62 +         grouped_rewards_with_diversity = grouped_rewards +
grouped_diversity_bonuses
63
64 +         # Compute group statistics
65 -         group_means = grouped_rewards.mean(dim=1, keepdim=True)
66 +         group_means = grouped_rewards_with_diversity.mean(dim=1,
keepdim=True)
67 +         if normalize_by_std:
68 -             group_stds = grouped_rewards.std(dim=1, keepdim=True) +
advantage_eps
69 -             advantages = (grouped_rewards - group_means) /
group_stds
70 +             group_stds = grouped_rewards_with_diversity.std(dim=1,
keepdim=True) + advantage_eps
71 +             advantages = (grouped_rewards_with_diversity -
group_means) / group_stds
72 +         else:
73 -             advantages = grouped_rewards - group_means
74 +             advantages = grouped_rewards_with_diversity -
group_means
75
76 +         # Flatten advantages back to original shape
77 +         advantages = advantages.view(-1)
78 +
79 +         # Update raw rewards to include diversity bonuses for
metadata
80 +         raw_rewards_with_diversity = raw_rewards + diversity_bonuses
81
82 +         # Compute metadata statistics
83 +         metadata = {
84 -             "mean_reward": raw_rewards.mean().item(),
85 -             "std_reward": raw_rewards.std().item(),
86 -             "max_reward": raw_rewards.max().item(),
87 -             "min_reward": raw_rewards.min().item(),

```

```

88 +         "mean_reward": raw_rewards_with_diversity.mean().item(),
89 +         "std_reward": raw_rewards_with_diversity.std().item(),
90 +         "max_reward": raw_rewards_with_diversity.max().item(),
91 +         "min_reward": raw_rewards_with_diversity.min().item(),
92 +         "mean_advantage": advantages.mean().item(),
93 +         "std_advantage": advantages.std().item(),
94 +         "mean_diversity_bonus": diversity_bonuses.mean().item(),
95     }
96
97 -     return advantages, raw_rewards, metadata
98 +     return advantages, raw_rewards_with_diversity, metadata
99
100 def compute_naive_policy_gradient_loss(
101 --- repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/grpo.
    py
102 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/grpo.
    py
103 @@ -203,6 +203,6 @@
104     parser.add_argument("--eval_dataset_path", type=str, default
    = "../MATH/test.jsonl")
105     parser.add_argument("--output_dir", type=str, default="ckpts
    /")
106 -     parser.add_argument("--learning_rate", type=float, default=1
    e-5)
107 +     parser.add_argument("--learning_rate", type=float, default=3
    e-5)
108     parser.add_argument("--grpo_steps", type=int, default=200)
109     parser.add_argument("--group_size", type=int, default=8)
110     parser.add_argument("--rollout_subset_size", type=int,
    default=256)
111 @@ -212,7 +212,7 @@ if __name__ == "__main__":
112     parser.add_argument("--gradient_accumulation_steps", type=
    int, default=16)
113     parser.add_argument("--batch_size", type=int, default=4)
114     parser.add_argument("--cliprange", type=float, default=0.2)
115 -     parser.add_argument("--loss_type", type=str, default="
    grpo_clip")
116 +     parser.add_argument("--loss_type", type=str, default="
    reinforce_with_baseline")
117     parser.add_argument("--wandb_project", type=str, default="
    grpo-math")
118     parser.add_argument("--wandb_name", type=str, default="
    grpo_clip_1")
119     args = parser.parse_args()
120 --- repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/
    run_job.sh
121 +++ repo_variants_full_pipeline_bsz50_fixed2_epoch3/idea_22/
    run_job.sh
122 @@ -21,7 +21,7 @@ timeout 2h uv run \
123     --index https://download.pytorch.org/whl/cul28 \
124     --index-strategy unsafe-best-match \
125     python grpo.py \
126 -         --learning_rate 1e-5 \
127 +         --learning_rate 3e-5 \
128         --grpo_steps 20 \
129         --group_size 8 \
130         --rollout_subset_size 128 \
131 @@ -30,7 +30,7 @@ timeout 2h uv run \
132     --gradient_accumulation_steps 16 \
133     --batch_size 4 \
134     --cliprange 0.2 \
135 -     --loss_type grpo_clip \
136 +     --loss_type reinforce_with_baseline \

```

---

```
137     --wandb_name $wandb_name
138
139  echo "Experiment finished successfully!"
```