
Effective Reinforcement Learning for Reasoning in Language Models

Lianghuan Huang^{* 1 2} Shuo Li^{* 1} Sagnik Anupam¹ Insup Lee¹ Osbert Bastani¹

Abstract

Reinforcement learning (RL) has emerged as a promising strategy for improving the reasoning capabilities of language models (LMs) in domains such as mathematics and coding. However, most modern RL algorithms were designed to target robotics applications, which differ significantly from LM reasoning. We analyze RL algorithm design decisions for LM reasoning, for both accuracy and computational efficiency, focusing on relatively small models due to computational constraints. Our findings are: (i) on-policy RL significantly outperforms supervised fine-tuning (SFT), (ii) PPO-based off-policy updates increase accuracy instead of reduce variance, and (iii) removing KL divergence can lead to concise generations and higher accuracy. Furthermore, we find that a key bottleneck to computational efficiency is that the optimal batch sizes for inference and backpropagation are different. We propose a novel algorithm, DASH, that performs *preemptive sampling* (i.e., sample a large batch and accumulate gradient updates in small increments), and *gradient filtering* (i.e., drop samples with small advantage estimates). We show that DASH reduces training time by 83% compared to a standard implementation of GRPO without sacrificing accuracy. Our findings provide valuable insights on designing effective RL algorithms for LM reasoning.¹

1. Introduction

Recent advancements have shown that reinforcement learning (RL) algorithms can significantly enhance the mathematical reasoning capabilities of language models (LMs) (DeepSeek-AI et al., 2025a; Qwen et al., 2025; Zeng et al., 2025). Despite these results, there has been little systematic understanding of how different RL design decisions contribute to their effectiveness in the LM reasoning setting. Many of these algorithms were originally designed for robotics, while LM reasoning exhibits qualitatively different learning patterns, meaning different design decisions may be more effective (Ahmadian et al., 2024); indeed, even the space of relevant design decisions may be different for LM reasoning compared to robotics. Our goal is to answer the following question: **How do we design effective RL algorithms for improving the reasoning capabilities of LMs?** Importantly, we are interested not only in the performance (i.e., the final accuracy), but also efficiency (i.e., how quickly the algorithm converges). Furthermore, we focus on relatively small models (0.5B, 1.5B, and 3B) where we can explore a variety of different RL algorithms.

We perform a systematic analysis of the different design decisions in an RL algorithm. We start by considering the two most prevalent types of algorithms: supervised fine-tuning (SFT) (Chen et al., 2023; Zeng et al., 2023), also known as behavior cloning, and on-policy RL (e.g., policy gradient (Sutton et al., 1999), PPO (Schulman et al., 2017a), GRPO (Shao et al., 2024), etc.). While SFT is much more efficient, we find it to be significantly less effective at improving reasoning ability for the models we consider; this may be due to the inability for smaller models to effectively mimic the reasoning traces of larger models or humans. In contrast, we find that on-policy RL is highly effective at improving performance.

Next, we compare different kinds of on-policy RL algorithms. Compared to the original policy gradient (PG) algorithm, PPO is designed to improve stability by “freezing” the inference policy and taking multiple gradient steps. We find that while PPO can improve accuracy, it has significantly higher variance compared to PG, which is the opposite of conventional wisdom. PPO also introduces a KL divergence term to regularize the training policy; however, we find that it leads to

^{*}Equal contribution ¹Department of Computer Science, University of Pennsylvania, Philadelphia (PA), US ²Department of Physics and Astronomy, University of Pennsylvania. Correspondence to: Osbert Bastani <obastani@seas.upenn.edu>.

Accepted at *Methods and Opportunities at Small Scale (MOSS)*, ICML 2025, Vancouver, Canada.

¹Our code is here: https://github.com/shuoli90/efficient_reasoning.

lengthier generations and worse performance.

While on-policy RL is highly effective, existing algorithms are computationally expensive to run. Analyzing their performance bottlenecks, we find that the sampling procedure is a key bottleneck. The key issue is that inference and training require significantly different batch sizes to make optimal use of computational resources. Thus, it is much more effective to perform inference in a single large batch, and then accumulate gradient steps for this batch over multiple training steps. This strategy allows us to perform efficient sampling while using the PG algorithm. Combined with strategies to filter out samples with small advantage estimates, we call the resulting algorithm Distributed-Aggregated Sampling Handler (DASH). Compared to GRPO, DASH reduces on-policy training time by **83%** without sacrificing accuracy (Figure 1). We open-source DASH to facilitate further research.

To summarize, our key findings are as follows: 1, For the models we consider, we find on-policy RL to be effective but not SFT (Section A.2 and 2.1); 2, We propose DASH, which accelerates on-policy training by **83%** without compromising accuracy (Section A.3, 2.4 and 2.5); 3, We find that while PPO-style gradient updates can slightly improve accuracy, it can introduce instability into training (Section A.4 and 2.2); 4, We find that removing KL divergence regularization can lead to more concise generations and higher accuracies (Section A.5).

2. Effective RL for LM Reasoning

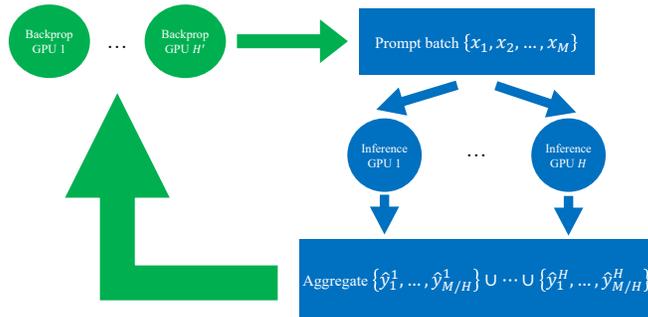


Figure 2. Illustration of preemptive sampling. We use H GPUs for inference and H' for backpropagation; they are shown in blue and green, respectively. Given a batch of M prompts $\{x_1, \dots, x_M\}$. The inference GPUs then generate corresponding responses $\{\hat{y}_1, \dots, \hat{y}_M\}$, which are aggregated across GPUs into CPU memory. When a backpropagation GPU requests generations for a prompt x_m , the corresponding cached response y_m is retrieved and delivered. Since we are using groups for advantage estimation, each prompt x_m is duplicated to form groups, and all generations in the same group are sent to the backpropagation GPU upon request.

First, we describe basic design decisions of our RL algorithm rooted in the prior literature; these are based either on experiments from prior work or our own experiments. Specifically, we consider an LM π_θ with parameters θ , which takes in a user prompt x and generates a reasoning trace \hat{y} , which we call a *trajectory*. We let \hat{y}_t denote the t th token in trajectory \hat{y} . For a training prompt x_n , we can check whether a generated trajectory \hat{y}_n produces the correct answer, represented as a scalar reward $r_n = R(x_n, \hat{y}_n) \in \mathbb{R}$. We assume that r_n is for the entire trajectory; typically, it is a binary indicator of

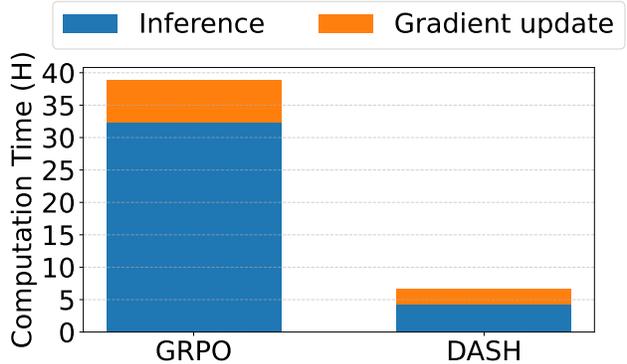


Figure 1. DASH can reduce running time by 83% compared to GRPO by using *preemptive sampling* (Section 2.4) and *gradient filtering* (Section 2.5).

whether the final answer is correct.²

2.1. RL Strategy

The first decision is what kind of RL strategy to use. We consider two strategies: supervised finetuning (SFT) and on-policy RL. SFT is effectively the same as behavior cloning, a popular imitation learning algorithm. Given a prompt training set $D = \{\mathbf{x}_n\}_{n=1}^N$, SFT collects corresponding expert trajectories $\mathcal{Y} = \{\mathbf{y}_n\}_{n=1}^N$, either from a human or a stronger LM. Then, the LM is optimized via maximizing the log likelihood on (D, \mathcal{Y}) :

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \log \pi_{\theta}(\mathbf{y}_n | \mathbf{x}_n) \quad \pi_{\theta}(\mathbf{y}_n | \mathbf{x}_n) = \prod_{t=1}^T \pi_{\theta}(\hat{y}_{n,t} | \mathbf{x}_n, \hat{y}_{n,1}, \dots, \hat{y}_{n,t-1})$$

Alternatively, on-policy RL learns from trajectories generated by the current LM π_{θ} . Given the a prompt set D , a typical on-policy RL algorithm optimizes the expected reward:

$$\pi_{\theta^*} = \arg \max_{\theta} J(\theta) \quad J(\theta) = \frac{1}{N} \sum_{\mathbf{x}_n \in D} \mathbb{E}_{\hat{\mathbf{y}}_n \sim \pi_{\theta}(\cdot | \mathbf{x}_n)} [R(\mathbf{x}_n, \hat{\mathbf{y}}_n)].$$

While SFT has been shown to be effective in settings such as Muennighoff et al. (2025), they use a post-trained larger sized model (Qwen2.5-32B-Instruct); our experiments show that it can be ineffective when the gap between the expert and π_{θ} is too large (specifically, we use small base models instead of larger instruction-tuned models). For instance, an expert may take leaps of reasoning that are incomprehensible to the learner. Thus, DASH uses on-policy RL. Another alternative that has been studied in the literature is self-imitation (Oh et al., 2018), where ‘‘expert’’ trajectories are obtained by performing search guided by π_{θ} , but results applying this strategy to LMs have so far been mixed (Shao et al., 2024).

2.2. Gradient Update Strategy

Next, we discuss the gradient update strategy. We consider both policy gradient (PG) (Sutton et al., 1999) and PPO (Schulman et al., 2017b) (which includes GRPO (Shao et al., 2024)). In general, we consider gradient approximations $\nabla_{\theta} J(\theta) \approx N^{-1} \sum_{n=1}^N J_n$ where J_n encodes the gradient approximation for the n th summand of $J(\theta)$. First, by the Policy Gradient Theorem, using

$$J_n^{\text{PG}} = \mathbb{E}_{\mathbf{y}_n \sim \pi_{\theta}(\cdot | \mathbf{x}_n)} \left[\frac{\nabla_{\theta} \pi_{\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n)}{\pi_{\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n)} A^{\pi_{\theta}}(\mathbf{x}_n, \hat{\mathbf{y}}_n) \right] \quad (1)$$

is exact, i.e., $\nabla_{\theta} J(\theta) = N^{-1} \sum_{n=1}^N J_n^{\text{PG}}$. Here, $A^{\pi_{\theta}}(\mathbf{x}_n, \hat{\mathbf{y}}_n)$ is the *advantage function*, which we discuss below. This update is truly on-policy since the trajectories $\hat{\mathbf{y}}$ must be sampled using the current policy π_{θ} . In robotics, PG can be unstable due to high variance when estimating the gradient $\nabla_{\theta} \pi_{\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n)$; as a consequence, π_{θ} can change rapidly across gradient steps, sometimes even becoming worse. PPO was devised to mitigate this instability. Specifically, they weaken the on-policy requirement, and ‘‘freeze’’ the data-generating policy $\pi_{\theta_{\text{old}}}$ for some number of gradient steps. The resulting update has the alternative form

$$J_n^{\text{PPO}} = \mathbb{E}_{\hat{\mathbf{y}}_n \sim \pi_{\theta_{\text{old}}}(\cdot | \mathbf{x}_n)} \left[\frac{\nabla_{\theta} \pi_{\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n)}{\pi_{\theta_{\text{old}}}(\hat{\mathbf{y}}_n | \mathbf{x}_n)} A^{\pi_{\theta_{\text{old}}}}(\mathbf{x}_n, \hat{\mathbf{y}}_n) \right],$$

where the differences compared to (1) are highlighted in red. Because this gradient is only valid when $\theta \approx \theta_{\text{old}}$, a KL regularization is imposed, to obtain $J_n^{\text{PPO-KL}} = J_n^{\text{PPO}} + \beta J_n^{\text{KL}}$, where $J_n^{\text{KL}} = \nabla_{\theta} D_{\text{KL}}(\pi_{\theta_{\text{base}}}(\cdot | \mathbf{x}_n) \| \pi_{\theta}(\cdot | \mathbf{x}_n))$. Following Jaques et al. (2019); Ouyang et al. (2022), the KL divergence term is with respect to the original model $\pi_{\theta_{\text{base}}}$ instead of $\pi_{\theta_{\text{old}}}$ as in PPO.

Critically, in PPO, θ_{old} is updated to be θ every K steps, where K is a hyperparameter. To further improve stability, the gradient is often clipped. GRPO uses the same gradient update as PPO; early versions include a weight $1/\text{len}(\hat{\mathbf{y}}_n)$ on the n th term to normalize by the length of the trajectory, but this term was removed in later versions (Liu et al., 2025). Finally, we note that when $\theta = \theta_{\text{old}}$, this gradient update is equivalent to the PG update (1); this property holds even with gradient clipping.

²Recent work has found that process rewards (Wang et al., 2024) may not be effective in our setting due to the difficulty predicting whether a reasoning trace is on the right track (DeepSeek-AI et al., 2025a).

Now, assume we have sampled a batch of M samples $\{(\mathbf{x}_m, \hat{\mathbf{y}}_m)\}_{m=1}^M$ from $\pi_{\theta_{\text{old}}}$, where initially $\theta = \theta_{\text{old}}$. If we take a single gradient step on all examples, then PPO coincides with PG. This is the strategy used by DASH. We consider two implementations of PPO that do not devolve into PG. First, we can take K gradient steps using all M samples, which we call *PPO-Multi* (or just *Multi*). Second, we can divide the M examples into K mini-batches of size M/K each, and take one gradient step on each mini-batch, which we call *PPO-Mini* (or just *Mini*). In our experiments, we find that DASH is more stable than Multi and Mini, suggesting that the added complexity of PPO-based off-policy gradient updates increases variance.

We include the KL term in DASH (i.e., $J_n^{\text{DASH}} = J_n^{\text{PG}} + \beta J_n^{\text{KL}}$) for closer comparison of DASH to Multi and Mini. However, in our experiments, we find that omitting the KL term improves accuracy.

2.3. Advantage Estimation

A key challenge in RL is estimating the quantity $A^\pi(\hat{\mathbf{y}} | \mathbf{x})$, which is called the *advantage* (Sutton & Barto, 2018); it is defined to be $A^\pi(\hat{\mathbf{y}} | \mathbf{x}) = Q^\pi(\hat{\mathbf{y}} | \mathbf{x}) - V^\pi(\mathbf{x})$, where Q^π is the Q-function and V^π is the value function. Intuitively, it captures how the specific generation $\hat{\mathbf{y}}$ compares to a random sample $\hat{\mathbf{y}}' \sim \pi_\theta(\cdot | \mathbf{x})$. In general, A^π is not known and must be estimated from data. We consider three strategies: (i) training a model to predict A^π , (ii) a Monte Carlo estimate called the *single-path method*, and (iii) a Monte Carlo estimate introduced by GRPO. The first approach is to train a model to predict $Q^\pi(\hat{\mathbf{y}} | \mathbf{x})$, which can be used to compute V^π and A^π (Schulman et al., 2017a). This approach can reduce variance, but recent work has found that it is highly biased due to the difficulty in predicting Q^π for reasoning tasks (Liang et al., 2022). Thus, we focus on Monte Carlo approaches.

The most popular Monte Carlo approach is the *single-path method*, which uses the estimate $A^{\pi_\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n) \approx r_n - \frac{1}{N} \sum_{n'=1}^N r_{n'}$, i.e., it is the centered reward; $b = N^{-1} \sum_{n'=1}^N r_{n'}$ is called the *baseline*. Intuitively, r_n is an estimate of the Q-function, and b is an estimate of the value function. A standard modification is to normalize by the standard deviation; this normalization can be useful when rewards tend to increase significantly as learning progresses, but our rewards are bounded so this cannot happen. Another modification is to leave out the reward for rollout n when estimating the value for rollout n , which reduces bias (Sutton & Barto, 2018); this modification can be important when N is small (e.g., $N = 2$) but only has a minor impact for larger N since the bias is small.

A shortcoming of the single-path method is that b is an estimate of the average value $N^{-1} \sum_{n'=1}^N V(\mathbf{x}_{n'})$ across all samples, whereas it ideally should estimate the value $V(\mathbf{x}_n)$. One alternative is the *vine method* (Kazemnejad et al., 2024; Schulman et al., 2017a), which uses a targeted sampling strategy to fix this issue; however, the vine method requires a large number of samples, making it computationally expensive. GRPO uses an advantage estimate that interpolates between the single-path and vine methods. It exploits the fact that in the reasoning setting, we typically train on multiple samples $\hat{\mathbf{y}}_n$ for a single user prompt \mathbf{x}_n . In our formulation, we can think of there being multiple \mathbf{x}_n that are identical. Suppose that we partition N into groups N_1, \dots, N_K , where \mathbf{x}_n is the same for all $n \in N_k$. Then, it estimates the advantage using the formula

$$A^{\pi_\theta}(\hat{\mathbf{y}}_n | \mathbf{x}_n) \approx r_n - \frac{1}{N_k} \sum_{n' \in N_k} r_{n'}, \quad (2)$$

where N_k is the group containing n . In other words, it replaces the baseline with a state-dependent baseline $b(\mathbf{x}_n) = N_k^{-1} \sum_{n' \in N_k} r_{n'}$; now, $b(\mathbf{x}_n)$ is an unbiased estimate of $V(\mathbf{x}_n)$. This strategy can be viewed as performing a vine estimate of the advantage at state \mathbf{x}_n , but not at any other state. DASH uses the GRPO advantage estimate (Section 2).

2.4. Preemptive Sampling

A key feature of RL for LMs is that inference typically occurs on specialized inference servers such as vLLM (Kwon et al., 2023a). Importantly, inference is typically much more memory efficient than backpropagation, meaning much larger batches are optimal for inference compared to backpropagation. Empirically, sampling takes up a much larger portion of training time than backpropagation if performed in small batches (Figure 1). Thus, we propose *preemptive sampling*, where we sample a large number of trajectories in one batch, and then perform backpropagation on these samples in smaller batches. Preemptive sampling can be further sped up by using multiple inference servers in parallel (Figure 2). In practice, our method can be used for both on-policy and off-policy sampling, depending on algorithmic design choices, as detailed in Section 2.2. Figure 2 illustrates preemptive sampling. DASH uses preemptive sampling.

2.5. Gradient Filtering

Finally, we propose to drop examples with small advantage estimates (which is equivalent to clipping small advantage values to zero, effectively dropping them from the gradient update). If the advantage estimate is small, then the contribution to the gradient is likely to be small (unless $\nabla_{\theta}\pi_{\theta}(\hat{y}_n | \mathbf{x}_n)$ happens to be very large, which we find to be unlikely in practice). Intuitively, these are examples where the model either almost always gets the answer right (in which case there is nothing new to learn) or almost always gets it wrong (in which case the problem is currently too difficult to learn). In addition, even if we only drop advantages that are exactly zero, this strategy can provide a speedup since backpropagation still takes time to compute the gradients $\nabla_{\theta}\pi_{\theta}(\hat{y}_n | \mathbf{x}_n)$ before they are multiplied by $A^{\pi_{\theta}}(\hat{y}_n | \mathbf{x}_n) = 0$. DASH uses gradient filtering.

3. Conclusion

We have performed a careful theoretical and empirical analysis (Appendix A) of some key design decisions in RL algorithms for improving language model reasoning, focusing on computationally constrained scenarios. We believe that systematizing the study of RL for language model reasoning is key to designing more effective RL algorithms in this domain, which differs significantly from robotics targeted by existing RL algorithms such as PPO. Our study is a first step in this direction.

References

- Ahmadian, A., Cremer, C., Gallé, M., Fadaee, M., Kreutzer, J., Pietquin, O., Üstün, A., and Hooker, S. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms, 2024. URL <https://arxiv.org/abs/2402.14740>.
- Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., and Sutton, C. Program synthesis with large language models, 2021. URL <https://arxiv.org/abs/2108.07732>.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., and Hoeffler, T. Graph of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690, March 2024. ISSN 2159-5399. doi: 10.1609/aaai.v38i16.29720. URL <http://dx.doi.org/10.1609/aaai.v38i16.29720>.
- Chen, B., Shu, C., Shareghi, E., Collier, N., Narasimhan, K., and Yao, S. Fireact: Toward language agent fine-tuning, 2023. URL <https://arxiv.org/abs/2310.05915>.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S., Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L., An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X. Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K., Wang, Y. Q., Wei, Y. X.,

- Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Xu, Y., Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z., and Zhang, Z. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025a. URL <https://arxiv.org/abs/2501.12948>.
- DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Guo, D., Yang, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Zhang, H., Ding, H., Xin, H., Gao, H., Li, H., Qu, H., Cai, J. L., Liang, J., Guo, J., Ni, J., Li, J., Wang, J., Chen, J., Chen, J., Yuan, J., Qiu, J., Li, J., Song, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Xu, L., Xia, L., Zhao, L., Wang, L., Zhang, L., Li, M., Wang, M., Zhang, M., Zhang, M., Tang, M., Li, M., Tian, N., Huang, P., Wang, P., Zhang, P., Wang, Q., Zhu, Q., Chen, Q., Du, Q., Chen, R. J., Jin, R. L., Ge, R., Zhang, R., Pan, R., Wang, R., Xu, R., Zhang, R., Chen, R., Li, S. S., Lu, S., Zhou, S., Chen, S., Wu, S., Ye, S., Ye, S., Ma, S., Wang, S., Zhou, S., Yu, S., Zhou, S., Pan, S., Wang, T., Yun, T., Pei, T., Sun, T., Xiao, W. L., Zeng, W., Zhao, W., An, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Li, X. Q., Jin, X., Wang, X., Bi, X., Liu, X., Wang, X., Shen, X., Chen, X., Zhang, X., Chen, X., Nie, X., Sun, X., Wang, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yu, X., Song, X., Shan, X., Zhou, X., Yang, X., Li, X., Su, X., Lin, X., Li, Y. K., Wang, Y. Q., Wei, Y. X., Zhu, Y. X., Zhang, Y., Xu, Y., Xu, Y., Huang, Y., Li, Y., Zhao, Y., Sun, Y., Li, Y., Wang, Y., Yu, Y., Zheng, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Tang, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Wu, Y., Ou, Y., Zhu, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Zha, Y., Xiong, Y., Ma, Y., Yan, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Wu, Z. F., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Huang, Z., Zhang, Z., Xie, Z., Zhang, Z., Hao, Z., Gou, Z., Ma, Z., Yan, Z., Shao, Z., Xu, Z., Wu, Z., Zhang, Z., Li, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Gao, Z., and Pan, Z. Deepseek-v3 technical report, 2025b. URL <https://arxiv.org/abs/2412.19437>.
- Glazer, E., Erdil, E., Besiroglu, T., Chicharro, D., Chen, E., Gunning, A., Olsson, C. F., Denain, J.-S., Ho, A., de Oliveira Santos, E., Järvinemi, O., Barnett, M., Sandler, R., Vrzala, M., Sevilla, J., Ren, Q., Pratt, E., Levine, L., Barkley, G., Stewart, N., Grechuk, B., Grechuk, T., Enugandla, S. V., and Wildon, M. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai, 2024. URL <https://arxiv.org/abs/2411.04872>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset, 2021. URL <https://arxiv.org/abs/2103.03874>.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- Jaques, N., Ghandeharioun, A., Shen, J. H., Ferguson, C., Lapedriza, A., Jones, N., Gu, S., and Picard, R. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog, 2019. URL <https://arxiv.org/abs/1907.00456>.
- Kazemnejad, A., Aghajohari, M., Portelance, E., Sordoni, A., Reddy, S., Courville, A., and Roux, N. L. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment, 2024. URL <https://arxiv.org/abs/2410.01679>.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023a.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J. E., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention, 2023b. URL <https://arxiv.org/abs/2309.06180>.
- Liang, L., Xu, Y., McAleer, S., Hu, D., Ihler, A., Abbeel, P., and Fox, R. Reducing variance in temporal-difference value estimation via ensemble of deep networks, 2022. URL <https://arxiv.org/abs/2209.07670>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step, 2023. URL <https://arxiv.org/abs/2305.20050>.

- Liu, J., Xia, C. S., Wang, Y., and Zhang, L. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=lqvX610Cu7>.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding rl-zero-like training: A critical perspective, 2025. URL <https://arxiv.org/abs/2503.20783>.
- Muennighoff, N., Yang, Z., Shi, W., Li, X. L., Fei-Fei, L., Hajishirzi, H., Zettlemoyer, L., Liang, P., Candès, E., and Hashimoto, T. s1: Simple test-time scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- Oh, J., Guo, Y., Singh, S., and Lee, H. Self-imitation learning, 2018. URL <https://arxiv.org/abs/1806.05635>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Puri, R., Kung, D. S., Janssen, G., Zhang, W., Domeniconi, G., Zolotov, V., Dolby, J., Chen, J., Choudhury, M., Decker, L., Thost, V., Buratti, L., Pujar, S., Ramji, S., Finkler, U., Malaika, S., and Reiss, F. Codenet: A large-scale ai for code dataset for learning a diversity of coding tasks, 2021. URL <https://arxiv.org/abs/2105.12655>.
- Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., and Qiu, Z. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. Zero: Memory optimizations toward training trillion parameter models, 2020. URL <https://arxiv.org/abs/1910.02054>.
- Schulman, J., Levine, S., Moritz, P., Jordan, M. I., and Abbeel, P. Trust region policy optimization, 2017a. URL <https://arxiv.org/abs/1502.05477>.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms, 2017b. URL <https://arxiv.org/abs/1707.06347>.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.11366>.
- Shypula, A., Li, S., Zhang, B., Padmakumar, V., Yin, K., and Bastani, O. Evaluating the diversity and quality of llm generated content, 2025. URL <https://arxiv.org/abs/2504.12522>.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In Solla, S., Leen, T., and Müller, K. (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Wang, P., Li, L., Shao, Z., Xu, R. X., Dai, D., Li, Y., Chen, D., Wu, Y., and Sui, Z. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024. URL <https://arxiv.org/abs/2312.08935>.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- West, P. and Potts, C. Base models beat aligned models at randomness and creativity, 2025. URL <https://arxiv.org/abs/2505.00047>.

- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models, 2023a. URL <https://arxiv.org/abs/2305.10601>.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models, 2023b. URL <https://arxiv.org/abs/2210.03629>.
- Zeng, A., Liu, M., Lu, R., Wang, B., Liu, X., Dong, Y., and Tang, J. Agenttuning: Enabling generalized agent abilities for llms, 2023. URL <https://arxiv.org/abs/2310.12823>.
- Zeng, W., Huang, Y., Liu, Q., Liu, W., He, K., Ma, Z., and He, J. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild, 2025. URL <https://arxiv.org/abs/2503.18892>.

A. Experimental Results

We perform experiments showing that (i) on-policy RL significantly outperforms SFT (Section A.2), (ii) DASH significantly reduces running time compared to standard GRPO (Section A.3), (iii) PG gradient updates outperform PPO-based gradient updates (Section A.4), and (iv) removing KL divergence can lead to more concise more generations and higher accuracies (Section A.5).

A.1. Experimental Setup

We use Qwen2.5- $\{0.5B, 1.5B, 3B\}$ models as our base models, all of which are not post-trained (i.e., no instruction tuning). We use the MATH dataset (Hendrycks et al., 2021), with the MATH-500 split (Lightman et al., 2023), which contains 12,000 examples for training and 500 examples for evaluation. We additionally use the GSM8K dataset (Cobbe et al., 2021) for out-of-distribution evaluation, which contains 1,319 examples. Finally, we also perform some experiments in the coding domain using the MBPP+ dataset (Liu et al., 2023), a 378-problem subset of verified problems from the MBPP dataset (Austin et al., 2021); we use 264 problems for training and 114 for evaluation. Additional details are provided in Table 7.

A.2. SFT vs. On-Policy RL

We compare three algorithms: (i) SFT with human-written reasoning traces, denoted *SFT-H*, (ii) SFT with reasoning traces from Qwen2.5-7B-Instruct, denoted *SFT-M*, and (iii) DASH. Results are shown in Table 1. As can be seen, DASH improves performance both in-distribution and out-of-distribution, demonstrating that on-policy algorithms can efficiently learn mathematical reasoning skills that generalize across datasets. On the other hand, neither SFT-H nor SFT-M improve performance, with SFT-H significantly degrading both in-distribution and out-of-distribution performance. Intuitively, the substantial performance degradation caused by SFT-H can be attributed to the fact that human reasoning often omits many intermediate steps, which is especially problematic for smaller LMs.

For coding, we train on human programs in MBPP+. Results are shown in Table 2. As can be seen, DASH outperforms SFT in most cases, demonstrating the the general effectiveness of on-policy RL at improving the reasoning capabilities of LMs. To the best of our knowledge, these are among the first results to show that on-policy RL can improve code generation for smaller LMs.

Method	Size (B)	MATH (%)	GSM8K (%)
Base	0.5	22.6	30.3
	1.5	48.0	58.8
	3.0	58.8	66.0
SFT-H	0.5	8.0	7.2
	1.5	17.2	32.8
	3.0	24.0	30.6
SFT-M	0.5	24.0	22.7
	1.5	46.2	46.0
	3.0	53.0	66.0
DASH	0.5	27.2	31.1
	1.5	54.0	58.8
	3.0	64.6	64.6

Table 1. Comparison of SFT to on-policy RL on math.

Method	Size (B)	pass1 (%)	pass@8 (%)
BASE	0.5	2.6	22.8
	1.5	7.1	60.5
SFT-H	0.5	8.77	29.0
	1.5	19.3	42.1
DASH	0.5	11.4	40.4
	1.5	23.7	63.2

Table 2. Comparison of SFT to on-policy RL on coding.

A.3. DASH vs. GRPO

Next, we compare DASH to GRPO both in terms of accuracy and running time by training Qwen2.5-0.5B using both GRPO and DASH. We also use an ablation of DASH without gradient filtering, denoted *No-GF*. Results are shown in Table 3 and illustrated in Figure 1. As can be seen, DASH significantly reduces GRPO training time (from 39 hours to 6.6 hours) without

any significant reduction in performance, highlighting the effectiveness of preemptive sampling and gradient filtering.

Compared to No-GF, DASH reduces running time by 4% without any significant reduction in performance. The effectiveness of gradient filtering can be improved; see Appendix C. We additionally show results for coding in Table 4. For coding, we again see a significant speedup, though it is smaller since the generation length is much smaller so the gap in optimal inference and backpropagation batch sizes is smaller.

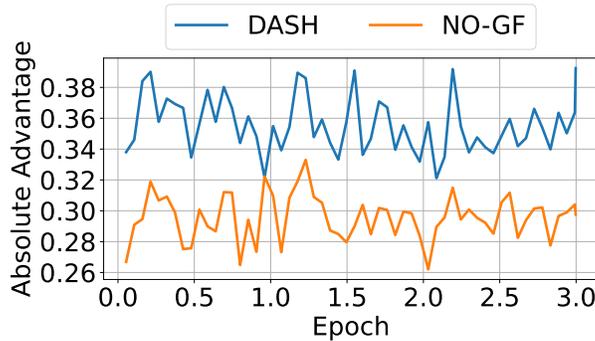
The impact of GF on training dynamics is illustrated in Figure 3. Specifically, as shown in Figure 3(a), gradient filtering increases the average absolute advantage values, leading to more significant gradient updates; consequently, as shown in Figure 3(b), forward and backward pass running times are reduced. Finally, since only samples inducing trivial gradient updates are filtered out, the training curves remain similar before and after applying gradient filtering, as shown in Figure 3(c).

Method	Time (h)	MATH (%)	GSM8K (%)
BASE	N/A	22.6	30.3
GRPO	38.9	27.6	32.8
No-GF	6.9	27.4	31.6
DASH	6.6	27.2	31.1

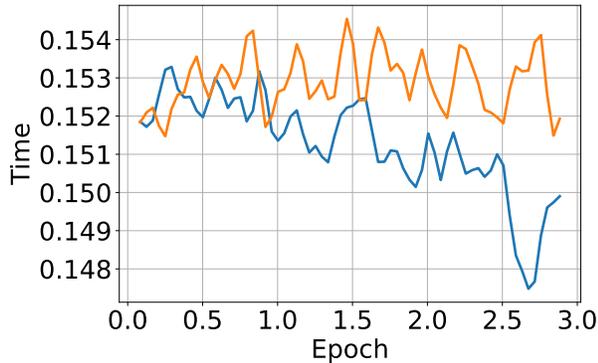
Table 3. Comparing on-policy RL algorithms on Qwen2.5-0.5B for math.

Method	Time (m)	pass@1 (%)	pass@8 (%)
BASE	N/A	2.3	22.8
GRPO	35.3	11.4	49.1
No-GF	16.3	10.5	43.9
DASH	16.5	11.4	40.4

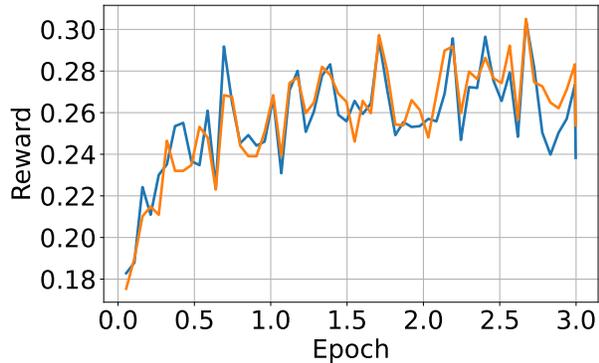
Table 4. Comparing on-policy RL algorithms using Qwen2.5-0.5B for coding.



(a) Average absolute advantage values.



(b) Loss computation time.



(c) Training reward recovers.

Figure 3. Comparison between DASH and No-GF for Qwen2.5-0.5B on math.

A.4. PG vs. PPO Gradient Updates

Next, we compare DASH to Multi and Mini. DASH uses a batch size of $M = 256$ (with $K = 1$), Multi uses $M = 256$ and $K = 3$, and Mini uses $M = 8$ so $K = 32$. Multi and Mini are slower than DASH; for fair comparison, we truncate their training times to match the wall-clock time of DASH. The results are shown in Table 5, and training curves are shown in Figure 4. As can be seen, Multi and Mini achieve faster initial performance improvements and have slightly higher accuracies; however, they have significantly more unstable training curves. Similar results for the 1.5B model are shown in Appendix C.

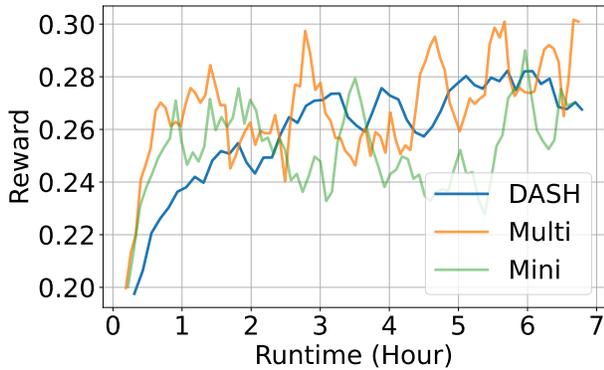


Figure 4. Training reward curves for PG vs. PPO on Qwen2.5-0.5B for math.

Method	Size (B)	MATH (%)	GSM8K (%)
DASH	0.5	27.2	31.1
	1.5	54.0	58.8
Multi	0.5	28.8	31.5
	1.5	54.0	61.0
Mini	0.5	29.8	31.6
	1.5	36.0	8.1

Table 5. Comparing PG to PPO for math.

A.5. KL Divergence Regularization

Next, we compare DASH to an ablation without the KL divergence term, denoted *No-KL*. Training reward curves are shown in Figure 6(a). As can be seen, removing KL divergence regularization generally leads to higher rewards during training; most likely, No-KL can focus on reward optimization without being constrained to stay close to the initial model. As shown in Table 6, No-KL achieves greater in- and out-of-distribution than DASH (except in the case of the out-of-distribution accuracy of the 3B model).

Furthermore, as shown in Figure 6(b), we find that for No-KL, the average generation length is shorter, thereby reducing overall training time; this difference is also reflected in Table 6. We hypothesize that to compensate for KL divergence regularization, models must generate longer reasoning traces.

Finally, for the 3B model, we study how KL divergence regularization affects pass@k. We follow Chen et al. (2021) to evaluate pass@k in an unbiased way. Results are in Figure 5: No-KL performs best for small k, although the gap closes for larger k. Intuitively, RL concentrates probability mass and reduces generation diversity (Shypula et al., 2025; West & Potts, 2025).

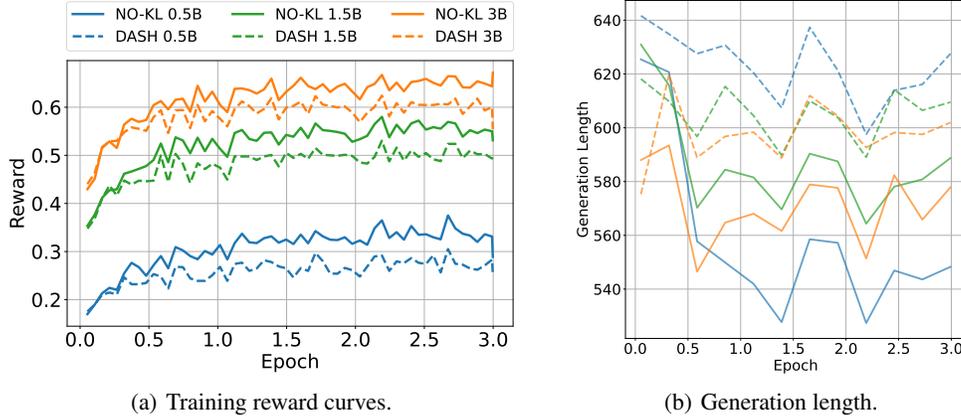


Figure 6. Comparing KL divergence regularization on math.

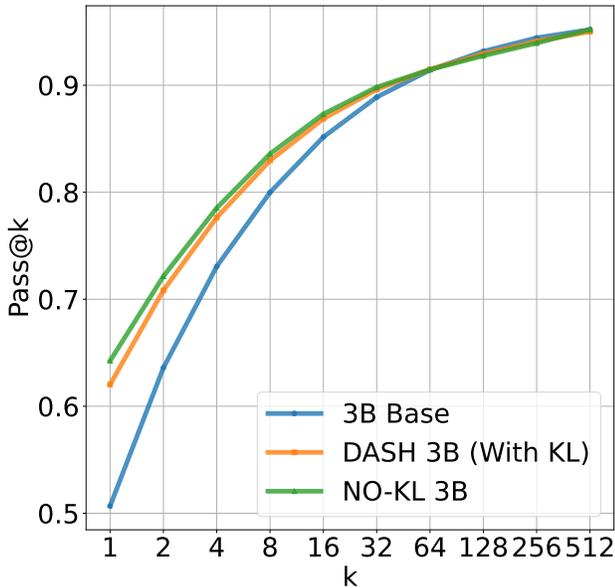


Figure 5. Impact of KL divergence regularization on pass@k for Qwen2.5-3B on math.

Method	Size (B)	Time (h)	MATH (%)	GSM8K (%)
Base	0.5	N/A	22.6	30.3
	1.5	N/A	48.0	58.8
	3.0	N/A	58.8	66.0
DASH (with KL)	0.5	6.6	27.2	31.1
	1.5	12.8	54.0	58.8
	3.0	22.6	64.6	64.6
No-KL	0.5	5.7	31.4	34.0
	1.5	10.3	56.8	62.1
	3.0	17.6	66.4	60.0

Table 6. Comparing No-KL to DASH and Base on math.

B. Additional Experimental Setup

Math. All GRPO experiments are conducted using 6 Nvidia A6000 GPUs; we use 4 GPUs for backpropagation and 2 for inference across all three model sizes (Qwen2.5- $\{0.5B, 1.5B, 3B\}$) (in practice, the 0.5B model only needs 2 GPUs for backpropagation, but we still use 4 for consistency).

Our implementation is based on Huggingface’s GRPO Trainer; the hyperparameters are as follows:

- Learning rate: 1e-06; for comparing to Multi and Mini on Qwen2.5-1.5B, we use 3e-06 due to larger batch size
- Backpropagation batch size per GPU: 2 (so batch size is 8)
- # generations per prompt: 4 (resulting in 2 prompts backpropagated on in each step)

- Maximum completion length: 2048
- Inference batch size for DASH: 256 (128 per inference GPU); for comparing to Multi and Mini on Qwen2.5-1.5B, we use 1024 (512 per inference GPU)
- Gradient accumulation steps for DASH: 32 (so the effective batch size is 256); for comparing to Multi and Mini on Qwen2.5-1.5B, we use 128 (so the effective batch size is 1024)
- Gradient steps per batch for Multi: 3
- Batch size for Mini: 8 (equivalently, no gradient accumulation)
- Gradient filtering threshold: 0.1

All other parameters are set to the default of the Huggingface trainer; a summary of the GRPO hyperparameters is in Table 7. To reduce memory footprint, we use DeepSpeed (Rajbhandari et al., 2020) ZeRO Stage 3 as well as CPU offload, gradient clipping, and mixed precision; our DeepSpeed configuration is shown in Figure 7.

Parameters for SFT are shown in Table 8. All SFT experiments use end-to-end fine-tuning instead of using parameter efficient methods such as LoRA (Hu et al., 2021). For model-generated reasoning traces, we use Qwen2.5-7B-Instruct as the teacher to keep the distribution of generations in the Qwen family. We use a temperature of 0.7 and filter out reasoning traces with the wrong answer. The resulting training set has 8,955 examples.

Versions of python and key libraries are shown in Table 9. The dev version of trl was cloned directly from trl’s GitHub repository on April 10, 2025.

Hyperparameter	Qwen2.5-0.5B	Qwen2.5-1.5B	Qwen2.5-3B
NVIDIA A6000 GPUs (training / sampling)	4 / 2 (2 / 2 using ZeRO)	4 / 2	4 / 2
Learning rate for DASH based runs	1×10^{-6}	1×10^{-6}	1×10^{-6}
Learning rate for Multi and Mini	1×10^{-6}	3×10^{-6}	N/A
Epochs	3	3	3
Batch size per device	2	2	2
Generations per prompt	4	4	4
Max completion length (tokens)	2048	2048	2048
Gradient accumulation steps for DASH based runs	32	32	32
Gradient accumulation steps for Multi and Mini	32	128	N/A
Gradient steps per sampled batch for Multi	3	3	N/A
Gradient-filtering threshold	0.1	0.1	0.1
Normalize gradients by generation length?	No	No	No

Table 7. Experimental configuration and hyperparameters for on-policy RL on MATH.

Parameter	Value	Package	Version
Learning rate	2×10^{-5}	python	3.11.11
Epochs	3	trl	0.17.0.dev0
Batch size per device	4	vllm	0.8.1
Gradient accumulation steps	2	pytorch	2.6.0

Table 8. Experimental configuration and hyperparameters for SFT.

Table 9. Package versions.

```

compute_environment: LOCAL_MACHINE
debug: false
deepspeed_config:
  gradient_clipping: 1.0
  offload_optimizer_device: cpu
  offload_param_device: cpu
  zero3_init_flag: false
  zero3_save_16bit_model: false
  zero_stage: 3
distributed_type: DEEPSPEED
downcast_bf16: 'no'
enable_cpu_affinity: false
machine_rank: 0
main_training_function: main
mixed_precision: bf16
num_machines: 1
num_processes: 4
rdzv_backend: static
same_network: true
tpu_env: []
tpu_use_cluster: false
tpu_use_sudo: false
use_cpu: false

```

Figure 7. DeepSpeed configuration.

Coding. All experiments with MBPP+ on coding using on-policy RL for Qwen2.5-0.5B were conducted on AWS EC2 g6.12xlarge instances with 48 vCPUs, 192 GiB memory, and 4 NVIDIA L4 Tensor Core GPUs with 96 GiB total GPU memory, with 2 GPUs dedicated to training and 2 to sampling. Experiments with Qwen2.5-1.5B and Qwen2.5-3B were conducted on AWS EC2 g6e.12xlarge instances with 48 vCPUs, 384 GiB memory, and 4 NVIDIA L40S Tensor Core GPUs with 192 GB total GPU memory, with 2 GPUs dedicated to training and 2 to sampling. All SFT experiments on MBPP+ were conducted using 2 NVIDIA A6000 GPUs. The hyperparameters for coding are the same as for math.

C. Additional Experiment Results

We compare gradient filtering with larger batch sizes, finding that gradient filtering is more effective when the per device batch size is 4 (instead of 2). This experiment is only possible for the for Qwen2.5-0.5B on the math dataset using our compute. Results are shown in Table 10 and training curves are shown in Figure 8. The time reduction achieved is larger than before (10% instead of 4%). These results suggest that gradient filtering may become more effective with larger batch sizes.

Method	Time (h)	MATH (%)	GSM8K(%)
No-GF	5.1	31.8	31.3
DASH	4.6	28.4	30.9

Table 10. Comparing No-GF to DASH with per device batch of 4 instead of 2 for Qwen2.5-0.5B on math.

We also show the comparison of Mini, Multi, and DASH on the 1.5B model (Figure 9). Conclusions are similar to Section A.4. In this case we stabilize Multi by increasing the gradient accumulation step to 128, but the high instability of Mini leads to decrease in training rewards as well as accuracies as shown in Table 5.

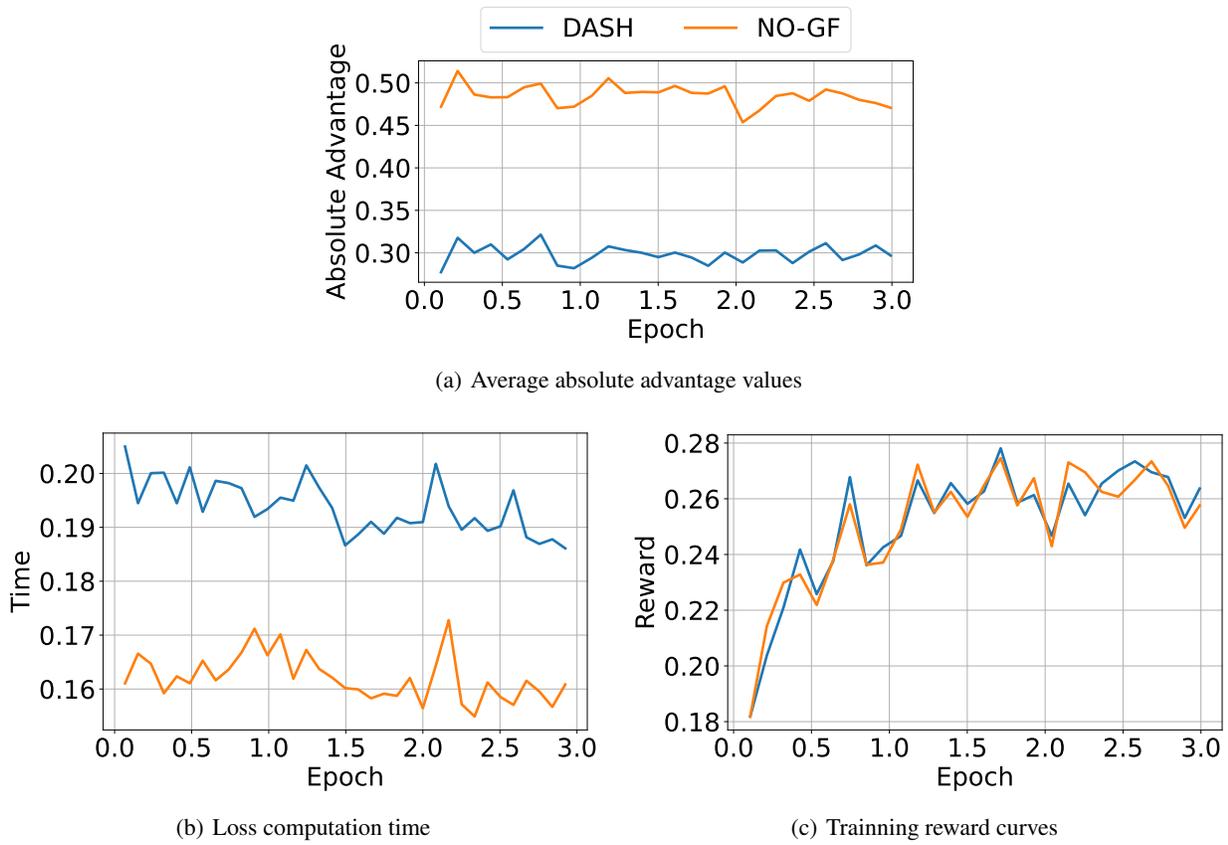


Figure 8. Comparison of No-GF and DASH with a per-device batch size of 4 for Qwen2.5-0.5B on math.

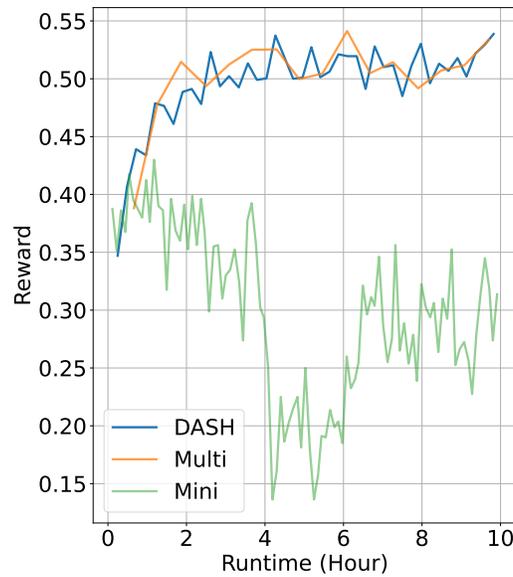


Figure 9. Training reward curves for PG vs. PPO on Qwen2.5-1.5B for math.

D. Related Work

LM Reasoning. Given the promising performance of language models (LMs), numerous studies have explored their application to mathematical problem solving (Hendrycks et al., 2021; Cobbe et al., 2021; Glazer et al., 2024), program synthesis (Austin et al., 2021; Puri et al., 2021), and other reasoning tasks. Since LMs often exhibit varying performance when directly prompted for these tasks, various methods have been proposed to explicitly elicit reasoning. For instance, Chain-of-Thought prompting (Wei et al., 2023) encourages LMs to generate intermediate reasoning steps before producing the final answer. Tree-of-Thought (Yao et al., 2023a) and Graph-of-Thought (Besta et al., 2024) extend this idea by imposing logical structure to organize the reasoning process. LM reasoning has also been enhanced through tool use (Yao et al., 2023b; Shinn et al., 2023). While these methods have proven effective in guiding LM reasoning and improving downstream task performance, they primarily focus on better prompt design rather than improving the models’ inherent reasoning capabilities.

RL for LM reasoning. Recent efforts have focused on using RL to improve LM reasoning capabilities. In question-answering tasks, FireAct (Chen et al., 2023) and AgentTuning (Zeng et al., 2023) enhance reasoning capabilities by learning from demonstrations from humans or stronger models. These approaches are commonly referred to as *supervised fine-tuning* (SFT), or *behavior cloning* in the RL literature. However, several studies have found limits on the effectiveness of SFT, instead proposing to use on-policy RL (DeepSeek-AI et al., 2025a; Shao et al., 2024; Zeng et al., 2025).

On the other hand, on-policy RL can be very computationally expensive, leading to a great deal of interest in improving efficiency. One shortcoming is that they require re-sampling generations after each model update, leading to sample inefficiency and prolonged training times. To mitigate this, DeepSeek-AI et al. (2025b) propose more efficient transformer architectures to accelerate pretraining, and Kwon et al. (2023b) introduce advanced memory management techniques to speed up sampling in post-training. Although current RL algorithms can leverage vLLM acceleration, the full potential of vLLM remains underutilized, leaving significant room for improving RL efficiency.