

Beyond Token-Level Policy Gradients for Complex Reasoning with Large Language Models

Anonymous ACL submission

Abstract

Existing policy-gradient methods for autoregressive language models typically select subsequent tokens one at a time as actions in the policy. While effective for many generation tasks, such an approach may not fully capture the structure of complex reasoning tasks, where a single semantic decision is often realized across multiple tokens—for example, when defining variables or composing equations. This introduces a potential mismatch between token-level optimization and the inherently block-level nature of reasoning in these settings. To bridge this gap, we propose **Multi-token Policy Gradient Optimization (MPO)**, a framework that treats sequences of K consecutive tokens as unified semantic actions. This block-level perspective enables our method to capture the compositional structure of reasoning trajectories and supports optimization over coherent, higher-level objectives. Experiments on mathematical reasoning and coding benchmarks show that MPO outperforms standard token-level policy gradient baselines, highlight the limitations of token-level policy gradients for complex reasoning, motivating future research to look beyond token-level granularity for reasoning-intensive language tasks.¹

1 Introduction

Large language models (LLMs) have become the foundation of modern natural language understanding and generation, achieving remarkable results through large-scale pretraining and autoregressive modeling (Kumar, 2024; Zhu et al., 2024; Tang et al., 2025). Recently, there has been increasing interest in leveraging policy gradient methods to further fine-tune these models, with the aim of enhancing their capacity for complex reasoning and long-horizon dependency modeling (Ouyang et al., 2022). Among these methods, Proximal

Policy Optimization (PPO) has emerged as a dominant framework, offering efficient and stable updates for reinforcement learning-based fine-tuning in LLMs (Schulman et al., 2017). Other advanced methods, such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024) and Decoupled Clip and Dynamic Sampling Policy Optimization (DAPO) (Yu et al., 2025), employ group-based and adaptive sampling strategies to address the challenges of structured reasoning and mathematical problem-solving more effectively.

However, as shown in Figure 1 left, in complex reasoning tasks such as mathematical reasoning, decision-making (like defining variables or completing equations) evolves over semantic units which contains multiple tokens. Existing policy gradient techniques decompose these structured reasoning steps into a series of local token predictions, disrupting the coherence of the underlying semantic actions (Figure 1 right). This fundamental granularity mismatch calls for optimization frameworks that transcend the limitations of token-level actions, enabling more semantically meaningful and globally coherent reasoning abilities in large language models (Mirzadeh et al., 2024; Zhang et al., 2025). Overcoming this issue calls for training strategies that capture reasoning actions at a coarser, semantically meaningful level beyond single-token updates.

In this paper, we introduce *Multi-token Policy Gradient Optimization (MPO)*, a framework that incorporates block-level action into the policy gradient process. As illustrated in Figure 1 right, instead of treating each token generation as an isolated action, MPO aggregates contiguous blocks of K tokens and unites importance sampling ratios over them. MPO allows the model to consider multiple correlated tokens as a single semantic action, better preserving the internal structure of reasoning steps such as variable definitions,

¹All codes and dataset will be released upon acceptance.

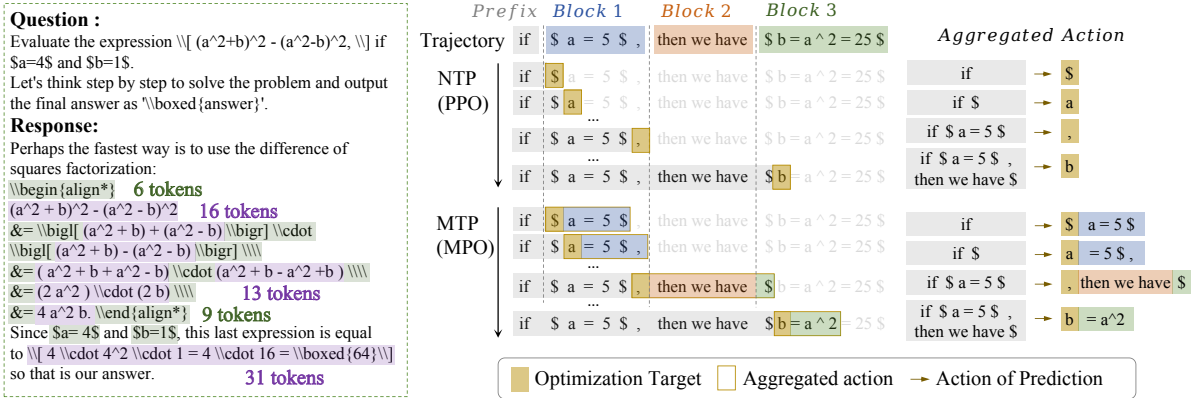


Figure 1: **Left:** In reasoning tasks such as mathematical problem-solving or code generation, the model’s decision process often spans across blocks of tokens—such as equations or functions—rather than being determined by each token independently; **Right:** illustration of token-level (NTP/PPO) vs. block-level (MTP/MPO) optimization. MPO aggregates K tokens as a semantically meaningful block for prediction and optimization, thereby better capturing sequence structure and long-range dependencies.

function calls, or equation formations. This block-level optimization encourages the policy to plan over meaningful reasoning segments rather than isolated symbols, thus maintaining consistency across intermediate decisions and improving global reasoning coherence. MPO only modifies the importance sampling ratio; it remains broadly compatible with existing policy-gradient frameworks and can be seamlessly integrated into contemporary LLM post-training pipelines. Our main contributions are:

1. We propose **Multi-token Policy Gradient Optimization (MPO)**, which enables the policy to optimize over structured reasoning units rather than isolated tokens.
2. We incorporate structural multi-token optimization into the post-training stage of LLMs, revealing its potential to enhance reasoning coherence and informing new directions for policy-gradient research.
3. MPO consistently outperforms token-level policy gradient baselines on mathematical reasoning and code generation benchmarks, demonstrating its effectiveness in improving reasoning ability during post-training.

2 Related Work

2.1 Multi-Token Prediction with LLMs

Multi-token prediction (MTP) is an extension of standard auto-regressive language modeling, where typically, the model is trained to predict only the

immediate next token given the preceding context. Gloeckle et al. (Gloeckle et al., 2024) propose using multiple prediction heads to forecast several tokens from each context position, improving sample efficiency and inference speed on coding and generative benchmarks. This approach demonstrates stronger induction and reasoning ability gains, particularly at larger model scales like DeepSeek V3 (Liu et al., 2024). Gerontopoulos et al. (Gerontopoulos et al., 2025) extend this to freeze pre-trained models and add learnable "register tokens" to support multi-token prediction without modifying the backbone.

While previous works on multi-token prediction enhanced LLM’s sample efficiency and accelerated inference, our work extends the MTP concept from pre-training or inference use into the realm of policy-gradient optimization.

2.2 Post-Training RL Algorithms

Reinforcement Learning from Human Feedback (RLHF) is an approach to align language models with human preferences by optimizing the policy using reinforcement learning on preference-labeled data. PPO (Schulman et al., 2017) introduced the clipped surrogate objective to stabilize updates, becoming the backbone of early RLHF such as InstructGPT (Ouyang et al., 2022). Building on PPO, approaches such as VinePPO (Kazemnejad et al., 2025) introduce intermediate rollouts to facilitate more accurate long-context credit assignment, thereby enhancing policy optimization

in autoregressive frameworks. Group Relative Policy Optimization (GRPO) (Shao et al., 2024) computes advantages by comparing multiple sampled completions per prompt (group-based advantage), eliminating the need for a value network and enhancing learning on reasoning tasks. DAPO (Decoupled Clip and Dynamic Sampling Policy Optimization) (Yu et al., 2025) further refines GRPO’s token-level loss by allowing asymmetric clipping bounds (“clip-higher”) and dynamic sampling to maintain a helpful gradient signal. VAPO (Yue et al., 2025) integrates value-based methods into reasoning and RLHF, augmenting PPO variants with adaptive GAE and value pretraining for enhanced stability and performance. GSPO (Zheng et al., 2025) (Group Sequence Policy Optimization) optimizes at the sequence level using importance sampling and clipping, improving stability and efficiency in policy optimization for LLMs. However, our MPO focuses on block-level actions, better aligning with the semantic reasoning units which are typically local blocks rather than entire sequences. CISPO (Chen et al., 2025) (Minimax-style policy optimization) methods have been proposed mainly in adversarial language alignment, thus falling outside our comparison.

While prior work enhanced PPO variants via sampling tweaks, decoupled clipping, or value modeling, to our knowledge, none have incorporated *multi-token action* in policy gradient optimization process as illustrated in Figure 1.

3 Preliminaries

We first review existing approaches to importance sampling in policy optimization—specifically PPO and its extensions such as GRPO—to clarify their formulation and highlight their limitations. This sets the stage for introducing our method, which extends importance sampling to account for longer-horizon behavior via multi-token predictions.

3.1 Importance Sampling Strategy in Policy Optimization

Proximal Policy Optimization (PPO) (Schulman et al., 2017) relies on importance sampling to support multiple epochs of minibatch updates with trajectories collected under a previous policy. The surrogate objective uses:

$$r_t = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}, \quad (1)$$

and the clipped PPO loss is

$$J(\theta) = \mathbb{E}_t \left[\min(r_t \hat{A}_t, \text{clip}(r_t, 1 \pm \epsilon) \hat{A}_t) \right], \quad (2)$$

where \hat{A}_t is the advantage estimate, and the clipping constrains large updates to maintain stability. Generalized Reward Policy Optimization (GRPO) (Shao et al., 2024) modifies the importance sampling ratio by introducing group-relative advantage sampling. The same importance ratio $r_t = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$ is used per sample, and the loss aggregates contributions averaged over sampled actions. Decoupled Clip and Dynamic Sampling Policy Optimization (DAPO) similarly builds on this framework by adapting the clipping threshold and employing rejection sampling, but does not alter the fundamental importance sampling ratio.

However, in the context of language modeling, where accurate semantics and the maintenance of long-range dependencies depend on structural token blocks, the commonly used token-level ratios fail to capture behavioral differences or dependencies that span multiple steps, and the standard per-step ratio may become highly volatile (Metelli et al., 2020; Papini et al., 2024).

3.2 Multi-Token Prediction Mechanism

To address the limited expressiveness of token-level representations, we turn to a multi-token representation mechanism. The MTP implementation introduced in DeepSeek-V3/R1 (Liu et al., 2024) allows the model at each position t to predict up to K tokens $o_{t+1}, o_{t+2}, \dots, o_{t+K}$ through a sequence of MTP modules that preserve causal consistency (we provide detailed implementation of MTP modules in Appendix A). Specifically, the probability distribution on position $t + k$ is:

$$h_t^k = L_k(M_k(h_t^{k-1}, \text{Emb}(o_{t+k-1}))), \quad (3)$$

where $h_t^k \in \mathbb{R}^h$ is the hidden state of the k^{th} MTP module, note that h_t^0 represents the hidden state output of the backbone model for token o_t . M_k is a learned projection layer, L_k is a Transformer decoder block, and $\text{Emb}(o_{t+k})$ is the embedding of token o_{t+k} . Each module outputs with an individual softmax prediction for the probability distribution of the $t + k + 1$ token of the sequence:

$$p(o_{t+k} | q, o_{1:t+k-1}) = \text{LM_Head}(h_t^k), \quad (4)$$

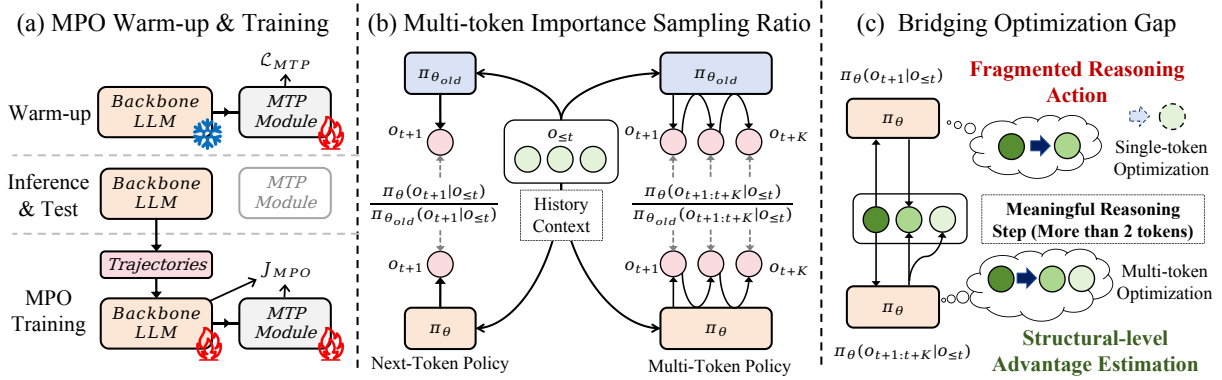


Figure 2: (a) Demonstration of the implementation of MPO warm-up and training process; (b) illustration of the united importance sampling ratio proposed in MPO method; (c) Comparison between single-token and multi-token optimization, multi-token optimization jointly models contiguous tokens as a structural reasoning action.

for $k = 2, \dots, K$. In supervised fine-tuning settings, the MTP objective is defined as:

$$\mathcal{L}_{\text{MTP}} = - \sum_{k=2}^K \alpha_k \log p(o_{t+k} | q, o_{1:t+k-1}). \quad (5)$$

The loss of the MTP module is assigned a decaying weight α_k to simulate the diminishing weights of tokens. Although multi-token prediction has been adopted in various pre-training and fine-tuning settings to improve generation efficiency and performance, its integration into reinforcement-based post-training remains limited. Existing policy-gradient approaches, including PPO, GRPO, and DAPO, remain confined to token-level updates and overlook the structured, multi-step action patterns inherent to reasoning-oriented generation.

Following previous works (Cai et al., 2024; Ankner et al., 2024), MPO first initializes the MTP modules using the last layer of the backbone model, then warm-up these modules using the objective (5) to ensure the quality of multi-token prediction, see Appendix A for more details.

4 Multi-token Objective

To address the limitations discussed in Section 3, we move beyond the conventional token-level optimization paradigm and introduce a framework that directly incorporates block-level semantic structure into policy gradient updates, thereby better aligning model optimization with the demands of complex reasoning.

4.1 From Single-Token to Multi-Token

Let q_i denote the prompt (input) for trajectory i , and $o_{i,1:t}$ denote the sequence of generated tokens

up to position t for trajectory i . The usual per-token importance sampling ratio $r_{i,t}$ for trajectory i at position t is defined as:

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t+1} | q_i, o_{i,1:t})}{\pi_{\theta_{\text{old}}}(o_{i,t+1} | q_i, o_{i,1:t})}. \quad (6)$$

With the help of multi-token prediction modules, we can replace the importance sampling ratio with an aggregated ratio over spans of length K (as shown in Figure 2(b)):

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t+1:t+K} | q_i, o_{i,1:t})}{\pi_{\theta_{\text{old}}}(o_{i,t+1:t+K} | q_i, o_{i,1:t})}. \quad (7)$$

The new objective then sums (or clips) these block-level ratios times advantages. This encourages coherent multi-token patterns rather than independent next-token moves. For example, when constructing the deepseek-r1 reasoning model, the pre-training stage cost around 90% of the total cost, and the post-training (supervised fine-tuning and then reinforcement learning alignment) cost the remaining 10% budget. For those pretrained LLMs that did not use MTP loss during their pre-training process, we provide a solution to update the model’s multi-token prediction ability in a relatively lower-cost post-training way.

4.2 Multi-Token Policy-Gradient Objective

Let each trajectory i produce a sequence of generated tokens using the original auto-regressive strategy $o_i = (o_{i,1}, o_{i,2}, \dots, o_{i,|o_i|})$. After generating the base sequence $o_{i,1:t}$ for prompt q_i , we compute the multi-token importance sampling ratio by sampling K additional tokens in an auto-regressive fashion. Specifically, at each position t , for each $n \in [1, K]$, the generation of token

$o_{i,t+n}$ is conditioned on the entire prefix $o_{i,1:t+n-1}$, meaning each new token is predicted based on all previously generated tokens, including those just sampled within the span. The revised multi-token importance sampling ratio is thus formulated as:

$$R_{i,t}^{(K)}(\theta) = \prod_{n=1}^K \frac{\pi_{\theta}(o_{i,t+n} | o_{i,1:t+n-1})}{\pi_{\theta_{\text{old}}}(o_{i,t+n} | o_{i,1:t+n-1})}. \quad (8)$$

Though equation 8 can effectively reduce bias by taking more future actions into consideration (see mathematical derivations in Appendix B.2), such a production form of the importance ratio tends to dramatically increase the variance of sampling ratios, which further leads to a large clip fraction and hinders the optimization process in practice. Inspired by the Log-COP-TD method (Hallak and Mannor, 2017), we adopt an alternative trade-off formulation to control variance, replacing the product of ratios with a weighted log-sum:

$$\begin{aligned} \tilde{R}_{i,t}^{(K)}(\theta) &= \exp\left(\sum_{n=1}^K \beta_n \log r_{i,t+n}(\theta)\right), \\ r_{i,t+n}(\theta) &= \frac{\pi_{\theta}(o_{i,t+n} | o_{i,1:t+n-1})}{\pi_{\theta_{\text{old}}}(o_{i,t+n} | o_{i,1:t+n-1})}, \end{aligned} \quad (9)$$

where β_n are non-negative step-wise weights satisfying $\sum_{n=1}^K \beta_n = 1$. Based on the weight of the first MTP module, denoted as β_2 , we define

$$\beta_k = \beta_2 \times \lambda^{k-2}, \quad k \geq 2; 0 \leq \lambda \leq 1, \quad (10)$$

where λ is a hyperparameter controlling the rate of information decay. This formulation applies a decaying weight to the $(k-1)^{\text{th}}$ MTP module, thereby incorporating the influence of multi-token information with diminishing strength. This design prioritizes the impact of nearer token predictions while still incorporating longer-horizon contributions in a controllable manner. We then propose the policy-gradient surrogate objective of MPO, analogous to PPO but using the K -step ratio at each position t :

$$\begin{aligned} J_{MPO}(\theta) &= \mathbb{E}_{q_i \sim D, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q_i)} \\ &\left[\frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min\left(\tilde{R}_{i,t}^{(K)}(\theta) \hat{A}_{i,t}, \right. \right. \\ &\left. \left. \text{clip}\left(\tilde{R}_{i,t}^{(K)}(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}\right) \hat{A}_{i,t}\right) \right], \end{aligned} \quad (11)$$

where $\hat{A}_{i,t}$ is the estimated advantage at position t . By jointly observing K consecutive tokens in a

single step, the estimator encourages optimization to account for a broader, structurally coherent view at the level of multi-token reasoning chunks (as illustrated in Figure 2(c)). In this paper, we focus on the experiments with $G = 1$ and we set $\epsilon_{\text{low}} = \epsilon_{\text{high}}$. We also conduct experiments with MTP-based value estimation in MPO, see Appendix C.

5 Experimental results

In this section, we provide a brief overview of our experimental setup, including the baselines, evaluation methods, and training hyperparameter configurations. And then, we analyze the task performance of MPO and the effects of introducing multi-token information and the training efficiency.

5.1 Experimental Settings

Model and Datasets. We implement MPO on three widely used instruction-tuned backbone models: Llama3.2-1B-Instruct (Dubey et al., 2024), DeepSeek-Distilled-Qwen2.5-1.5B, and DeepSeek-Distilled-Qwen2.5-7B (Guo et al., 2025), to evaluate its effectiveness across different architectures and model scales. We assess MPO on two mathematical reasoning benchmarks of varying difficulty, GSM8K (Cobbe et al., 2021) (college-level) and MATH (Hendrycks et al., 2021) (competition-level), and further evaluate its code generation capability on the HumanEval benchmark (Chen, 2021) to test its cross-domain generalization. For HumanEval, we use coding benchmark MBPP (Austin et al., 2021) as the training set. See Dataset statistics in Appendix A.2.

Evaluation. We use pass@1 metric to evaluate the models' accuracy in answering mathematical questions and coding completions, which focuses on the correctness of the final answer provided by the model upon completion of the reasoning process. For math problems, to unify the evaluation process, we use the answering format of "`\boxed{answer}`". We select the most commonly used token-wise policy optimization strategies: PPO (Schulman et al., 2017), GRPO (Shao et al., 2024), and DAPO (Yu et al., 2025). In this paper, we mainly discuss the implementation and results of MPO based on PPO.

Implementation. As mentioned in section 3.2, we adopt a cold-start setting for both our proposed method and baselines, only warming up the MTP module before the training process of MPO. See Appendix A for more implementation details.

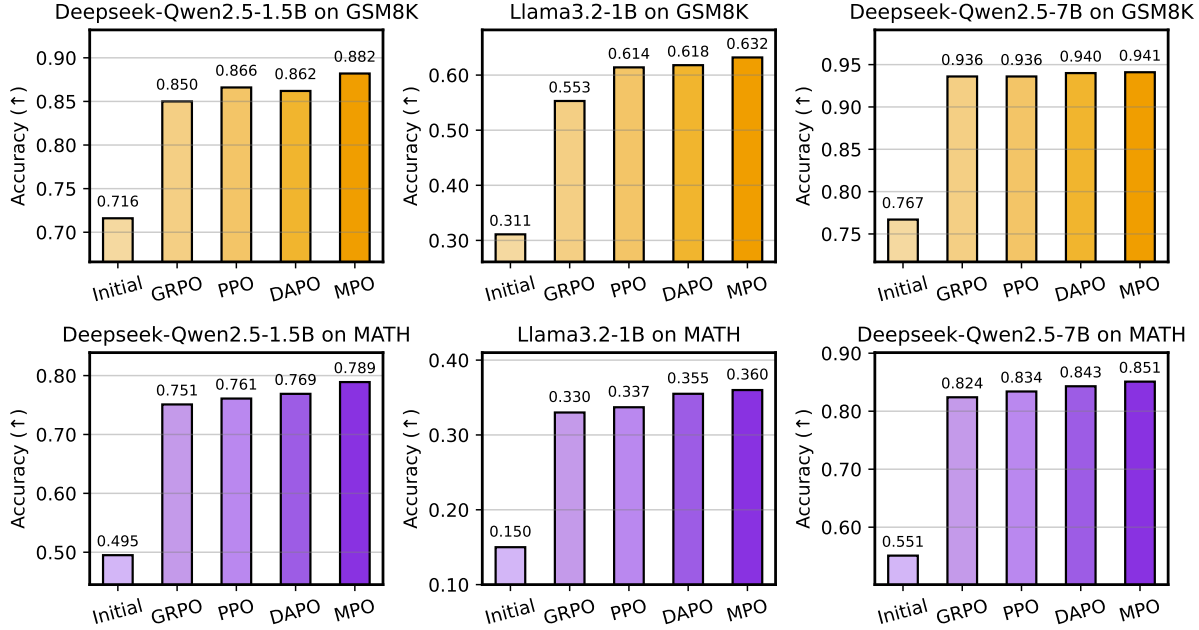


Figure 3: The performance of proposed MPO and baseline methods. MPO outperforms the baselines in most scenarios, demonstrating the effectiveness of aggregating block-wise information.

Model	Zero-Shot	GRPO	DAPO	PPO	MPO (Ours)
Llama3.2-1B-Instruct	0.354	0.372	0.396	0.390	0.403
DeepSeek-Qwen2.5-1.5B	0.451	0.591	0.603	0.598	0.640
DeepSeek-Qwen2.5-7B	0.689	0.811	0.817	0.805	0.841

Table 1: The performance of proposed MPO and baseline methods on coding task HumanEval, MPO consistently outperforms the baselines, demonstrating the effectiveness of aggregating block-wise actions during optimization.

5.2 Task Performance

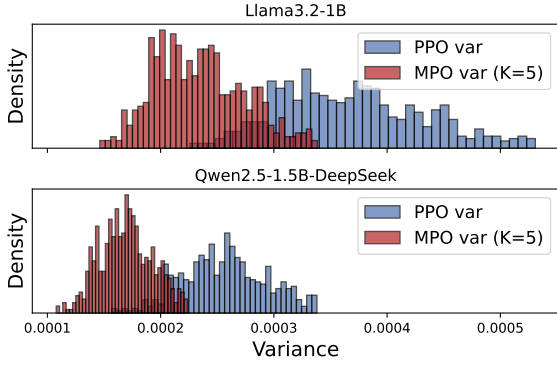
We conduct experiments to evaluate the performance of MPO. The results are shown in Figure 3 and Table 1. The proposed method consistently outperforms the baseline approaches—including PPO, GRPO, and DAPO—across both GSM8K and MATH benchmarks, under both evaluated model architectures and scales. This improvement demonstrates that optimizing over block-level semantic actions is beneficial for the model. On the HumanEval benchmark (as shown in Table 1), MPO also achieves steady gains over PPO, GRPO, and DAPO, confirming that its advantages are not limited to symbolic reasoning but also extend to program synthesis tasks. Therefore, in code generation tasks, considering multi-step semantic dependencies plays a crucial role in ensuring both accuracy and coherence in generation. Although the performance of MPO on GSM8K with a 7B model is close to the baseline, it achieves better results on more challenging benchmarks, such as MATH and HumanEval, indicating its advantage in structured reasoning tasks with larger models.

5.3 Analysis of Training Stability

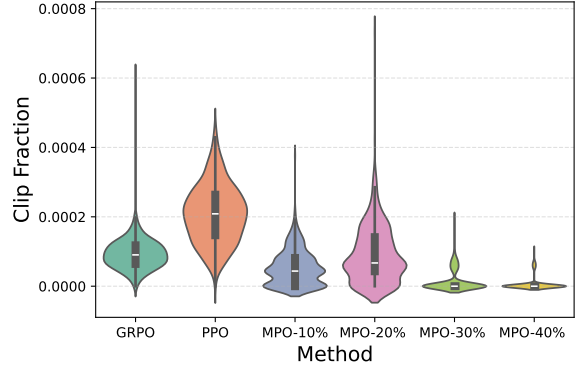
As discussed in Section 4.2, incorporating block-wise information broadens the optimization horizon of each action and helps reduce bias, particularly in mathematical reasoning and code generation tasks where semantic units naturally span multiple tokens. Empirically, we observe that MPO stabilizes the importance sampling ratio throughout training. As shown in Figure 4a and Figure 4b, analysis on GSM8K reveals that integrating multi-token information consistently lowers both the variance of the importance sampling ratio and the clip fraction during optimization. Notably, this stabilizing effect becomes more pronounced as the weights β_k associated with the MTP modules increase.

5.4 Ablation Study of MPO Weights

Another important question is how increasing the contribution of MTP weights influences learning dynamics. As shown in Figure 5, we examine the effect of varying the proportion of multi-token information incorporated via MTP modules on



(a) Variance of importance sampling ratio.



(b) IS Ratio Clip fraction among strategies.

Figure 4: Comparison of the variance of importance sampling ratios and clip fraction during training.

Dataset	MPO Task Accuracy		
	K=2	K=3	K=5
GSM8K	0.871	0.875	0.882
MATH	0.771	0.753	0.789

(a) Analysis of the number of MTP modules.

MTP Weight $\beta_2 =$	0.08	0.06	0.04
$\lambda = 1.0$	0.745	0.773	0.787
$\lambda = 0.9$	0.756	0.785	0.760
$\lambda = 0.8$	0.758	0.789	0.764

(b) Grid Search of beta and decay rate.

Table 2: Effect of MTP block size K and decay rate λ on training stability and performance. Extending the block size to $K = 5$ and applying a moderate decay $\lambda = 0.8$ produces the most stable and effective result.

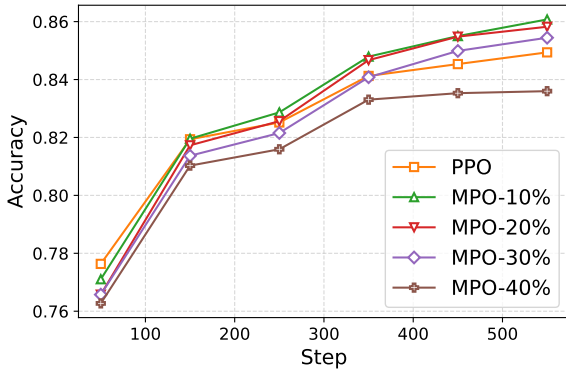


Figure 5: The effect of varying the proportion of weights incorporated from the MTP modules.

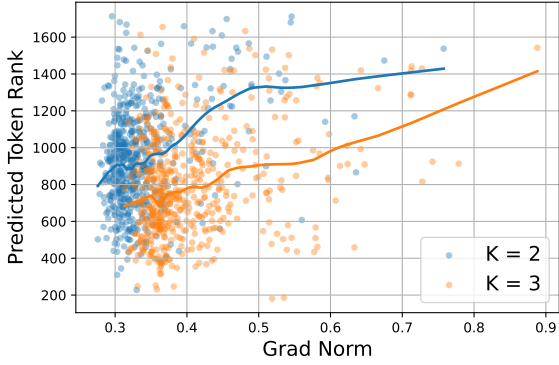
the GSM8K dataset using the DeepSeek-Distill-Qwen-1.5B model. The block-level contribution is controlled by adjusting the cumulative weights of the MTP modules (i.e., the sum of β_2 to β_K), while holding the decay factor fixed at $\lambda = 0.8$ as defined in Eq. 10. Increasing the proportion of future information initially enhances policy stability, as evidenced by reduced variance in importance sampling ratios and a lower frequency of gradient clipping. Optimal performance is achieved when approximately 10% 20% of the training signal originates from MTP modules. However, when this proportion becomes too large, the adopted ratio approximation and the noise introduced by multi-token prediction increase bias, ultimately weakening the alignment between gradient updates and the next-token action objective.

5.5 The Bias-variance Trade-off in MPO

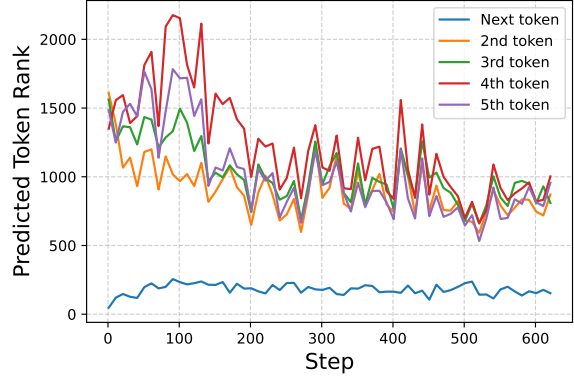
Having analyzed the stability of MPO under moderate integration of multi-token information, we conducted a grid search over the number of aggregated tokens K and the decay rate λ to assess their effects on training stability and optimization performance. As shown in Table 2a, increasing K up to 5 leads to improved performance, suggesting that extending the decision horizon within this range is beneficial; however, larger K values tend to introduce additional noise. Table 2b summarizes the effects of varying the decay rate λ and the cumulative MTP weight $\sum_{k=2}^K \beta_k$. Across all experiments, we keep the initial weight β_2 small so that the total block-level contribution remains below 30%, which stabilizes training while preserving meaningful multi-token context. If β_2 is too small (e.g., $\beta_2 = 0.04$), the stability of information from the MTP modules declines. In general, applying a decay term for more distant predictions improves stability, as appropriately discounting further predictions helps constrain variance while maintaining a sufficiently extended decision horizon—balancing bias reduction and training stability in MPO.

5.6 Noise Analysis of Multi-token Prediction

As discussed in section 5.4, excessive MTP weights were found to reduce overall performance. We further examine whether such instability



(a) Predicted token rank with avg. gradient norm.



(b) Predicted token rank with training steps.

Figure 6: Reliability analysis of the multi-token information. (a) Larger prediction errors correlate with higher gradient magnitudes. (b) Progressive improvement in the MTP module’s next-K-token prediction capability.

475 arises from noisy multi-token information. We
 476 use the logit rank of the predicted token—its
 477 position within the backbone model’s vocabulary
 478 distribution—as a proxy for prediction accuracy.
 479 Figure 6a shows the relationship between the
 480 2nd-token prediction rank and the average batch
 481 gradient norm under $K = 2$ and 3. Higher ranks
 482 correlate with larger gradient norms, suggesting
 483 that less accurate multi-token predictions can
 484 indeed amplify gradient noise and destabilize
 485 training. Figure 6b further illustrates that, although
 486 only warmed up initially, the MTP module’s
 487 predictions improve steadily during training even
 488 without explicit entropy regularization. The impact
 489 of prediction noise remains controlled.

5.7 Efficiency Analysis

491 **Training Efficiency.** Table 3 and Table 4
 492 summarize the resource overhead and efficiency
 493 of MPO on a single node with $8\times$ NVIDIA
 494 A100 (80 GB). For Table 4 we use Deepseek-
 495 Distill-Qwen2.5-1.5b as the base model, and
 496 we set group size=4 for GRPO and $K=5$ for
 497 MPO. The memory consumption of MPO is
 498 comparable to that of GRPO, while training runs
 499 approximately 30% faster on average. MPO
 500 contributes controlled training burden while
 501 offering substantial performance benefits. The
 502 additional parameters introduced by MTP are
 503 used solely during training to estimate multi-token
 504 policy ratios and are detached during inference,
 505 ensuring a fair comparison with baseline methods.
 506 **Warm-up Strategy.** In MPO, we include a brief
 507 warm-up phase; this design choice follows prior
 508 multi-head decoding frameworks(Cai et al., 2024;
 509 Ankner et al., 2024). The warm-up phase costs
 510 approximately 12 minutes while PPO or MPO

Training Speed (Second per iteration)			
Model	PPO	K=3	K=5
Deepseek-Qwen2.5-1.5b	25.6	30.3	33.3
Llama3.2-1b	19.1	25.0	28.3

Table 3: The average training time with block size K .

Method	Memory Usage	Time Per Step
PPO	22.1%	25.6 s
GRPO	30.6%	43.7 s
MPO	31.0%	33.3 s

Table 4: Memory and time efficiency.

511 training costs approximately 7 hours for 20 epochs
 512 on average. The additional time overhead is
 513 negligible compared to the main training stage.
 514 More discussion about efficiency is in Appendix D.

6 Conclusion

515 In this paper, we propose Multi-Token Policy
 516 Gradient Optimization (MPO), mitigating the
 517 gap between token-level optimization and the
 518 structured reasoning behavior required by complex
 519 reasoning tasks. Experiments on GSM8K, MATH,
 520 and HumanEval show that MPO consistently
 521 outperforms standard policy-gradient baselines.
 522 Further analyses demonstrate that moderate
 523 incorporation of multi-token information stabilizes
 524 training by reducing gradient variance, and that
 525 there exists an optimal integration ratio balancing
 526 bias variance trade-off. Our paper not only
 527 introduces a practical multi-token policy gradient
 528 algorithm, but also marks a methodological shift
 529 that moves beyond token-level optimization in the
 530 post-training of LLMs. We hope this study inspires
 531 future research to further explore the granularity of
 532 policy learning and advance toward semantically
 533 structured decision-making process in the LLMs.
 534

535 Limitations

536 MPO builds upon Multi-Token Prediction (MTP),
537 a technique primarily explored in pre-training
538 and speculative decoding. Its application to
539 post-training optimization remains understudied,
540 presenting several open research questions:

541 **1. Warmup Method.** While MTP-based training
542 has been shown to improve reasoning (DeepSeek
543 V3 and subsequent work), applying it during post-
544 training—especially to models not pre-trained with
545 MTP (e.g., Llama-3, Qwen2.5)—may degrade
546 final performance after reinforcement learning.
547 Although MPO mitigates this by freezing the
548 backbone, developing more effective warmup
549 strategies remains an open challenge.

550 **2. Scaling Block Size.** MPO incurs computational
551 overhead as the number of MTP modules increases.
552 Our experiments used up to $K=5$ modules, whereas
553 prior work suggests larger K (e.g., 8) may yield
554 further gains. Designing efficient methods to
555 leverage longer future contexts is crucial, especially
556 for scaling block-level policy optimization.

557 **3. Alternative Optimization Objectives.** This
558 work implements MPO based on a PPO-based
559 objective. As noted in Section 4.2, MPO could also
560 be integrated with other objectives (e.g., GRPO
561 and its variants), but implementing MPO based on
562 GRPOs may encounter challenges like rebalancing
563 the bias-variance trade-off. Investigating how
564 MPO interacts with different surrogate objectives
565 presents a promising direction for future work.

566 References

567 Zachary Ankner, Rishab Parthasarathy, Aniruddha
568 Nrusimha, Christopher Rinard, Jonathan Ragan-
569 Kelley, and William Brandon. 2024. Hydra:
570 Sequentially-dependent draft heads for medusa
571 decoding. *arXiv preprint arXiv:2402.05109*.

572 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten
573 Bosma, Henryk Michalewski, David Dohan, Ellen
574 Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1
575 others. 2021. Program synthesis with large language
576 models. *arXiv preprint arXiv:2108.07732*.

577 Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu
578 Peng, Jason D Lee, Deming Chen, and Tri Dao.
579 2024. Medusa: Simple llm inference acceleration
580 framework with multiple decoding heads. In
581 *International Conference on Machine Learning*,
582 pages 5209–5235. PMLR.

583 Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang,
584 Bo Fei, Bo Yang, Boji Shan, Changqing Yu, Chao
585 Wang, Cheng Zhu, and 1 others. 2025. Minimax-m1:

Scaling test-time compute efficiently with lightning
attention. *arXiv preprint arXiv:2506.13585*. 586 587

Mark Chen. 2021. Evaluating large language models
trained on code. *arXiv preprint arXiv:2107.03374*. 588 589

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, and 1 others. 2021. Training verifiers
to solve math word problems. *arXiv preprint
arXiv:2110.14168*. 590 591 592 593 594 595

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, and 1 others. 2024. The llama 3 herd of models.
arXiv e-prints, pages arXiv–2407. 596 597 598 599 600

Anastasios Gerontopoulos, Spyros Gidaris, and Nikos
Komodakis. 2025. Multi-token prediction needs
registers. *arXiv preprint arXiv:2505.10518*. 601 602 603

Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière,
David Lopez-Paz, and Gabriel Synnaeve. 2024.
Better & faster large language models via multi-token
prediction. In *Proceedings of the 41st International
Conference on Machine Learning*, pages 15706–
15734. 604 605 606 607 608 609

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong
Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
Deepseek-r1: Incentivizing reasoning capability in
llms via reinforcement learning. *arXiv preprint
arXiv:2501.12948*. 610 611 612 613 614 615

Assaf Hallak and Shie Mannor. 2017. Consistent on-
line off-policy evaluation. In *Proceedings of the
34th International Conference on Machine Learning-
Volume 70*, pages 1372–1383. 616 617 618 619

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021. Measuring mathematical
problem solving with the math dataset. *arXiv
preprint arXiv:2103.03874*. 620 621 622 623 624

Amirhossein Kazemnejad, Milad Aghajohari, Eva
Portelance, Alessandro Sordoni, Siva Reddy, Aaron
Courville, and Nicolas Le Roux. 2025. Vineppo:
Refining credit assignment in rl training of llms. In
*Forty-second International Conference on Machine
Learning*. 625 626 627 628 629 630

Pranjal Kumar. 2024. Large language models (llms):
survey, technical frameworks, and future challenges.
Artificial Intelligence Review, 57(10):260. 631 632 633

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,
Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi
Deng, Chenyu Zhang, Chong Ruan, and 1 others.
2024. Deepseek-v3 technical report. *arXiv preprint
arXiv:2412.19437*. 634 635 636 637 638

639	Alberto Maria Metelli, Matteo Papini, Nico Montali, and Marcello Restelli. 2020. Importance sampling techniques for policy optimization. <i>Journal of Machine Learning Research</i> , 21(141):1–75.	Yingjie Zhu, Xuefeng Bai, Kehai Chen, Yang Xiang, Jun Yu, and Min Zhang. 2024. Benchmarking and improving large vision-language models for fundamental visual graph understanding and reasoning. <i>arXiv preprint arXiv:2412.13540</i> .	695
640			696
641			697
642			698
643	Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. In <i>The Thirteenth International Conference on Learning Representations</i> .	A Implementation Details	700
644		We will provide separate explanations of the overall MPO implementation process, the experimental parameter settings, and details of the code and hardware implementation.	701
645			702
646			703
647			704
648			705
649	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	A.1 MTP Modules	706
650		The implementation of MPO requires the MTP module to predict the model’s policy distributions at multiple positions along a given trajectory i , enabling the calculation of importance sampling coefficients over blockwise sequences. Achieving stable multi-position decision making depends on the MTP module’s ability to produce accurate predictions of multi-token policy distributions. These distributions cannot be directly modeled using only the backbone’s final layer and language output head, so an appropriate warmup is necessary. Although previous works has noted that MTP can improve the model’s reasoning ability during fine-tuning (Liu et al., 2024), to ensure fair comparisons, we do not backpropagate gradients to the backbone during MPO finetuning and only train the MTP module (as shown in Figure 7), which has also been adopted in previous efficient MTP decoding works (Cai et al., 2024). As we discussed previously, each MTP module and its corresponding LM head are deeply copied from the backbone model’s final layer and LM head, respectively, to avoid an overly complex initialization process.	707
651			708
652			709
653			710
654			711
655			712
656	Matteo Papini, Giorgio Manganini, Alberto Maria Metelli, and Marcello Restelli. 2024. Policy gradient with active importance sampling. <i>arXiv preprint arXiv:2405.05630</i> .		713
657			714
658			715
659			716
660	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .		717
661			718
662			719
663			720
664	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .		721
665			722
666			723
667			724
668			725
669			726
670	Yihong Tang, Kehai Chen, Muyun Yang, Zhengyu Niu, Jing Li, Tiejun Zhao, and Min Zhang. 2025. Thinking in character: Advancing role-playing agents with role-aware reasoning. <i>arXiv preprint arXiv:2506.01748</i> .		727
671			728
672			729
673			730
674			731
675	Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. <i>arXiv preprint arXiv:2503.14476</i> .	A.2 Dataset and Training Parameters	732
676		We conduct experiments using three widely used mathematical reasoning and coding benchmarks: GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), and HumanEval (Chen, 2021). For HumanEval, we combined the sanitized and original data provided in coding benchmark MBPP (Austin et al., 2021) to conduct MPO training. Table 5 summarizes the main statistics of these datasets. The question-answer pairs in both datasets are provided in a natural language format, suitable for language model-based reasoning experiments. The open source mathematical and coding benchmarks used do not contain any private information or offensive content.	733
677			734
678			735
679			736
680	Yu Yue, Yufeng Yuan, Qiyong Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiase Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, and 1 others. 2025. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks. <i>arXiv preprint arXiv:2504.05118</i> .		737
681			738
682			739
683			740
684			741
685			742
686	Xiang Zhang, Juntao Cao, Jiaqi Wei, Yiwei Xu, and Chenyu You. 2025. Tokenization constraints in llms: A study of symbolic and arithmetic reasoning limits. <i>arXiv preprint arXiv:2505.14178</i> .		743
687			744
688			
689			
690	Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, and 1 others. 2025. Group sequence policy optimization. <i>arXiv preprint arXiv:2507.18071</i> .		
691			
692			
693			
694			

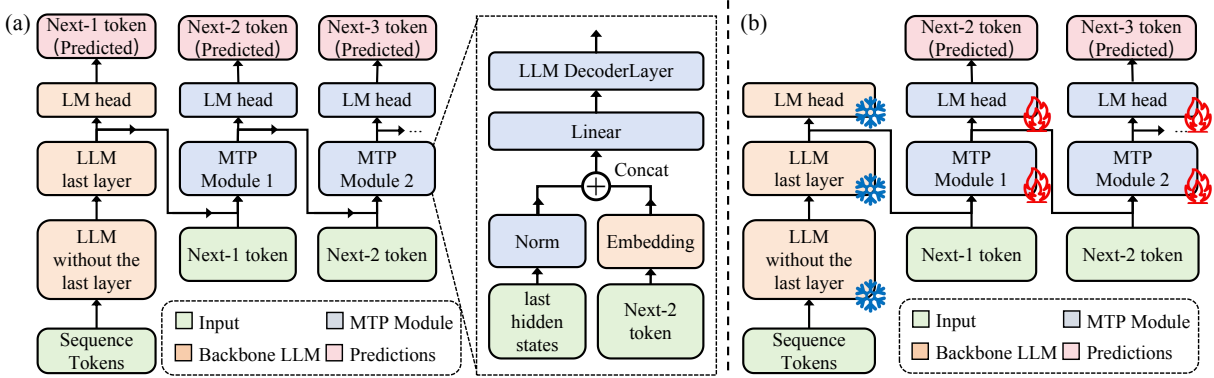


Figure 7: (a) The implementation of the MTP module used in MPO, the structure preserves causal consistency following Deepseek-V3; (b) Demonstration of the MPO warmup stage to initialize MTP Modules before the RL stage.

Dataset	# Train Samples	# Test Samples	Avg. Sample Length
GSM8K	7,473	1,319	~80 tokens
MATH	7,500	5,000	~200 tokens
HumanEval	1399 (MBPP)	164	~130 tokens

Table 5: Statistics of datasets used for training and evaluation.

We shuffle and concatenate the training splits of two math datasets for MTP warm-up training. For MPO training, we conduct separate experiments on each dataset, using the corresponding training split for learning and the corresponding test split for evaluation. This setup allows us to evaluate the generalization and effectiveness of our method on both mathematical reasoning and coding tasks.

A.3 MPO Training Process

Here we provide additional implementation details for the Multi-Token Policy Optimization (MPO) method described above, including the hyperparameters and a pseudocode Algorithm 1 for clarity. Table 6 lists our main hyperparameters and settings for the MPO training process. These values are extracted from our training scripts and reflect typical settings for large-scale RLHF or policy optimization tasks on math reasoning datasets. For HumanEval, due to the extensive thinking behavior of Deepseek distilled models, we extend the response length to 3072.

B Discussion of Multi-Token Approach

This appendix provides a formal discussion of why incorporating Multi-token prediction into policy optimization can improve optimization performance and enhance long-horizon reasoning ability. The exposition is consistent with

the notation and definitions in the main text (Sections 3.2–4).

B.1 Intuitive Example and Motivation

In standard policy optimization methods such as PPO, the policy improvement step is based on per-token decisions, where the importance ratio

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t+1} | q_i, o_{i,1:t})}{\pi_{\theta_{\text{old}}}(o_{i,t+1} | q_i, o_{i,1:t})} \quad (12)$$

compares the likelihood of a single action (next token) under the updated and old policies. This formulation implicitly assumes that: 1) the reward structure can be well captured by local token-level updates, and 2) long-horizon dependencies are either negligible or already encoded in the advantage estimator $\hat{A}_{i,t}$.

While traditional one-step token-level policy optimization updates the model based on individual token predictions, our illustration highlights that semantics and reasoning patterns in language generation typically emerge from coherent blocks of multiple tokens. As shown in the figure, generating a valid reasoning step or factual entity (e.g., \$a=5\$ or \$b = a^2 = 25\$) requires accurately predicting a sequence of tokens as a whole.

Limiting optimization to one-step advantage estimation introduces several issues:

Parameter	Value
Train batch size	256
PPO mini batch size	128
PPO micro batch size / GPU	4
Number of GPUs	8
Number of Nodes	1
Max prompt length	1024
Max response length	2048
Learning rate (actor)	1×10^{-6}
Learning rate (critic)	1×10^{-5}
Total epochs	20

Table 6: Key hyperparameters for MPO training.

Algorithm 1: Multi-Token Policy Optimization (MPO)

Input : Dataset \mathcal{D} of prompts q_i and annotated outputs o'_i . Pretrained LLM with MTP module; initial parameters θ . Old policy parameters θ_{old} . Block size K , weights α_k, β_k ; clip parameter $\epsilon_{\text{low}}, \epsilon_{\text{high}}$

Output : Updated policy parameters θ

Stage 1: MTP Warm-Up

for several epochs **do**

for batch of (q_i, o'_i) in \mathcal{D} **do**

for each position t in o'_i **do**

for $k = 1$ to K **do**

 Predict $p(o'_{t+k} | q, o'_{1:t})$ using backbone model ($k = 1$) and MTP module ($k \geq 2$)

 Compute $\mathcal{L}_{\text{MTP}} = -\sum_{k=2}^K \alpha_k \log p(o'_{t+k} | q, o'_{1:t})$

 Update only MTP parameters by minimizing \mathcal{L}_{MTP}

Stage 2: MPO Fine-tuning

for each MPO epoch **do**

$\mathcal{B} \leftarrow$ sample mini-batch of prompts q_i from \mathcal{D}

for each q_i in \mathcal{B} **do**

 Generate trajectory o_i using $\pi_{\theta_{\text{old}}}$

 Compute per-token or group-based advantages $\hat{A}_{i,t}$

for each trajectory (q_i, o_i) and each position t in o_i **do**

 // Sample K -step target block

 Sample $o_{i,t+1:t+K}$ in an autoregressive manner using $\pi_{\theta_{\text{old}}}$

 Initialize $D_{kl} = 0$

for $n = 1$ to K **do**

 Compute:

$p_{\theta} = \pi_{\theta}(o_{i,t+n} | o_{i,1:t+n-1})$

$p_{\theta_{\text{old}}} = \pi_{\theta_{\text{old}}}(o_{i,t+n} | o_{i,1:t+n-1})$

$D_{kl} += \beta_n \cdot (\log p_{\theta} - \log p_{\theta_{\text{old}}})$

$\tilde{R}_{i,t}^{(K)} = \exp(D_{kl})$

 Compute surrogate loss for (i, t) :

$J_{i,t} = \min \left(\tilde{R}_{i,t}^{(K)} \hat{A}_{i,t}, \text{clip}(\tilde{R}_{i,t}^{(K)}, 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}) \hat{A}_{i,t} \right)$

 Aggregate loss: $J_{\text{MPO}} = \frac{1}{N} \sum_i \sum_t J_{i,t}$.

 Update parameters θ via gradient descent using J_{MPO} . Optionally, update θ_{old} periodically.

- 797 1. **Fragmented learning signals:** Token-level
798 updates isolate local token actions, weak-
799 ening the internal consistency of reasoning
800 segments.
- 801 2. **Incomplete credit assignment:** Single-token
802 objectives struggle to capture dependencies
803 whose effects span across multiple tokens
804 within a reasoning step.
- 805 3. **Granularity mismatch:** The optimization
806 signal operates at the token level, while the
807 semantic correctness of reasoning typically
808 emerges at the segment level.

809 In contrast, as illustrated in the lower MTP/MPO
810 section, aggregating actions and optimization
811 targets at the block level enables updates that
812 better capture semantic integrity, reduce estimation
813 variance, and align training with the goals of
814 structured reasoning.

815 By instead defining a K -step decision ratio,

$$816 R_{i,t}^{(K)}(\theta) = \prod_{n=1}^K \frac{\pi_{\theta}(o_{i,t+n} \mid o_{i,1:t+n-1})}{\pi_{\theta_{\text{old}}}(o_{i,t+n} \mid o_{i,1:t+n-1})}, \quad (13)$$

817 we directly compare the joint likelihood of a
818 future span under the new and old policies. This
819 transforms the optimization from purely myopic
820 next-token correction into a multi-step decision
821 process, allowing the policy gradient to align with
822 behaviors that unfold across multiple tokens. In
823 essence, the surrogate objective is no longer tied
824 to local moves but reflects higher-level planning,
825 where corrections propagate over spans of length
826 K rather than being confined to single steps.

827 B.2 Theoretical Derivation and Bias 828 Reduction

829 B.2.1 Problem Setup and Notation

830 We consider a standard reinforcement learning
831 setup for sequence generation. Let time steps
832 be indexed by $t = 0, 1, 2, \dots$, where each step
833 corresponds to generating a token.

- 834 • State: s_t denotes the state at time t , which
835 includes the context (previous tokens and
836 hidden states).
- 837 • Action a_t denotes the token generated at time
838 t .
- 839 • Policy: $\pi_{\theta}(a_t \mid s_t)$, with π_{old} denoting the
840 previous policy.

- 841 • Discount factor: $\gamma \in (0, 1]$.

- 842 • Value function:

$$843 V^{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s \right]$$

- 844 • Advantage function:

$$845 A^{\pi}(s_t, a_t) = Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$$

- 846 • Estimated value function: $\hat{V}_{\phi}(s)$

847 In standard one-step PPO, the advantage is
848 often estimated using the TD error and generalized
849 advantage estimator (GAE):

$$850 \delta_t = r_t + \gamma \hat{V}_{\phi}(s_{t+1}) - \hat{V}_{\phi}(s_t), \quad (14)$$

$$851 \hat{A}_t^{\text{GAE}} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}, \quad (15)$$

852 where $\lambda \in [0, 1]$ controls the trade-off.

853 B.2.2 K-Step / Multi-Token Core Idea

854 In K -step (multi-token) prediction, at time t we
855 consider predicting not just a_t , but the next K
856 tokens jointly. Define the K -step return:

$$857 G_t^{(K)} := \sum_{k=0}^{K-1} \gamma^k r_{t+k} + \gamma^K \hat{V}_{\phi}(s_{t+K}). \quad (16)$$

858 The corresponding K -step advantage estimator is

$$859 \hat{A}_t^{(K)} := G_t^{(K)} - \hat{V}_{\phi}(s_t). \quad (17)$$

860 Intuitively, as K increases, the advantage relies
861 more on the actual observed rewards and less on the
862 bootstrapped value function \hat{V}_{ϕ} , thereby reducing
863 the bias from value function approximation.

864 B.2.3 Bias Analysis

865 Let $\epsilon_V := \sup_s \left| \hat{V}_{\phi}(s) - V^{\pi}(s) \right|$ denote the
866 maximum value function approximation error. The
867 K -step return based on the true value function is

$$868 G_t^{(K), \text{true}} := \sum_{k=0}^{K-1} \gamma^k r_{t+k} + \gamma^K V^{\pi}(s_{t+K}), \quad (18)$$

869 so the error due to \hat{V}_{ϕ} is

$$870 \Delta G_t^{(K)} := G_t^{(K)} - G_t^{(K), \text{true}} \\ 871 = \gamma^K (\hat{V}_{\phi}(s_{t+K}) - V^{\pi}(s_{t+K})). \quad (19)$$

872 Hence, in absolute value,

$$|\Delta G_t^{(K)}| \leq \gamma^K \epsilon_V. \quad (20)$$

The bias of the K-step advantage relative to the true K-step advantage is

$$\begin{aligned} \hat{A}_t^{(K)} - A_t^{\text{true},(K)} &= (G_t^{(K)} - \hat{V}_\phi(s_t)) \\ &\quad - (G_t^{(K),\text{true}} - V^\pi(s_t)) \\ &= \gamma^K (\hat{V}_\phi(s_{t+K}) - V^\pi(s_{t+K})) \\ &\quad - (\hat{V}_\phi(s_t) - V^\pi(s_t)), \end{aligned} \quad (21)$$

with absolute value bounded by

$$\left| \hat{A}_t^{(K)} - A_t^{\text{true},(K)} \right| \leq \gamma^K \epsilon_V + \epsilon_V = (1 + \gamma^K) \epsilon_V. \quad (22)$$

If we focus on the relative dependence on one-step ($n=1$), note that the error term for one-step is $(1 + \gamma)\epsilon_V$. More importantly, if we decompose the bias into “the bootstrapping bias from the tail” and “the bias from the current estimate,” we can see that the bootstrapping term from the tail is multiplied by γ^K .

In particular, when we consider the bias propagation with respect to the true infinite-horizon return, the bias term related to n -step bootstrapping decays as $\mathcal{O}(\gamma^K)$. In other words, the bias is introduced by bootstrapping (relying on $\hat{V}_\phi(s_{t+n})$), and its coefficient is γ^K . As n increases, this term decays exponentially by a factor of γ^K , thus reducing reliance on $\hat{V}_\phi \Rightarrow$ reducing the bias introduced by function approximation.

B.2.4 Expected Advantage Bias

Consider the expected bias over trajectories:

$$\begin{aligned} &\mathbb{E}[\hat{A}_t^{(K)} - A^\pi(s_t, a_t)] \\ &= \mathbb{E}\left[G_t^{(K)} - \hat{V}_\phi(s_t) - (Q^\pi(s_t, a_t) - V^\pi(s_t))\right] \\ &= \mathbb{E}\left[\sum_{k=0}^{K-1} \gamma^k r_{t+k} + \gamma^K \hat{V}_\phi(s_{t+K}) - \hat{V}_\phi(s_t) \right. \\ &\quad \left. - \sum_{k=0}^{\infty} \gamma^k r_{t+k} + V^\pi(s_t)\right] \\ &= \mathbb{E}\left[\gamma^K (\hat{V}_\phi(s_{t+K}) - V^\pi(s_{t+K})) \right. \\ &\quad \left. - (\hat{V}_\phi(s_t) - V^\pi(s_t)) \right. \\ &\quad \left. - \sum_{k=K}^{\infty} \gamma^k (r_{t+k} - \mathbb{E}[r_{t+k} | s_{t+K}])\right]. \end{aligned} \quad (23)$$

Under standard assumptions, the last term has zero expectation, so the dominant contribution to

bias comes from

$$\begin{aligned} &\mathbb{E}[\hat{A}_t^{(K)} - A^\pi(s_t, a_t)] \\ &\approx \mathbb{E}\left[\gamma^K (\hat{V}_\phi(s_{t+K}) - V^\pi(s_{t+K})) \right. \\ &\quad \left. - (\hat{V}_\phi(s_t) - V^\pi(s_t))\right]. \end{aligned} \quad (24)$$

If the bias of \hat{V}_ϕ at different time steps is not systematically correlated (or its expectation is approximately zero), then the main remaining error comes from the tail error scaled by γ^n , with the overall expected bias reduced to $\mathcal{O}(\gamma^n \epsilon_V)$. The expected bias due to reliance on the bootstrapped value in n -step is reduced to the order of γ^n , thereby lowering the bias.

C Multi-token Value Estimation

C.1 Multi-Token Value Heads

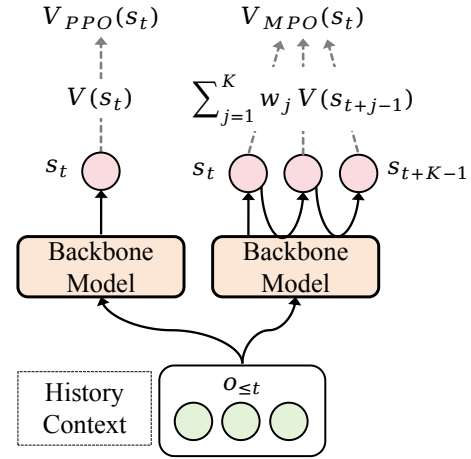


Figure 8: The implementation of the value function on the original auto-regressive backbone model (left) and those with the MTP module (right). For multi-token prediction models, we try a weighted sum of the output from the shared value head.

In parallel with our multi-token perspective modification for the actor, we have also conducted experiments on the critic model, allowing it to anticipate the returns of actions over a future interval during training. To accommodate multiple multi-token prediction heads, we aggregate their outputs at each decision point via a learned convex combination:

$$V_{\text{MPO}}(s_t) = \sum_{j=1}^K w_j V_j(\hat{s}_{t+j-1}) \quad (25)$$

where \hat{s}_{t+j-1} denotes the predicted states and

$$w_j = \frac{\exp(\gamma_j)}{\sum_{k=1}^K \exp(\gamma_k)}, \quad \gamma = (\gamma_1, \dots, \gamma_K) \quad (26)$$

Settings	GSM8K			MATH		
	K=2	K=3	K=5	K=2	K=3	K=5
MPO	0.871	0.875	0.882	0.771	0.753	0.789
+ V_{MPO}	0.879	0.875	0.882	0.773	0.756	0.762

Table 7: MPO performance with different token block size K and ablation study of multi-token value estimation.

with γ as learnable parameters, balancing the value obtained from predicted states. V_j is the value head append to the j^{th} MTP module. Given this value estimation, the advantage at each position t is computed as

$$\hat{A}_{i,t} = \hat{Q}_{i,t} - V_{\text{MPO}}(s_t), \quad (27)$$

where $\hat{Q}_{i,t}$ is the estimated return as in standard policy gradient methods. This approach permits each value head to specialize, while the aggregation provides a flexible baseline for stable advantage calculation and policy optimization over multi-token predictions. Figure 8 illustrates the value estimation process.

C.2 Effect of Block size on Value Estimation

We conducted experiments with Deepseek-Distill-Qwen2.5-1.5b by fixing the weight sum of MTP modules ($\sum_{k=2}^K \beta_k$) of knowledge at 20%, with results shown in Table 7. The method of injecting multi-token information from the critic side did not improve overall performance as expected. Although our initial experiments suggested that this approach could have a "symmetric" effect—namely, by training the MTP module on the value prediction task to enhance overall training stability—this was not supported by later results. After addressing certain implementation issues in our code, it appears that MPO is not sensitive to the incorporation of multi-token information into the value function. In fact, with $K = 4$ or 5 , it may introduce greater instability and cause the model to oscillate at a suboptimal level. Simply injecting value predictions by summing the outputs of multiple MTP modules does not seem to effectively improve value estimation accuracy.

D Discussion of Efficiency

Computation overhead. Table 8 illustrates further analysis with the resource overhead and comparative efficiency of MPO. The current implementation of MTP in MPO does not k-fold increase the overall computational cost. Since all

Training Speed (Second per iteration)			
Model	PPO	K=3	K=5
Deepseek-Qwen2.5-1.5b	25.6	30.3	33.3
Llama3.2-1b	19.1	25.0	28.3
Model Size (Compared to the original model)			
Model	PPO	K=3	K=5
Deepseek-Qwen2.5-1.5b	1.0×	1.20×	1.41×
Llama3.2-1b	1.0×	1.52×	2.05×

Table 8: The size of plug-and-play MTP modules and average training time per step.

logits needed for multi-token ratios are obtained from a single standard forward pass, the MTP module generates predictions for K tokens in parallel within one forward pass, without requiring K separate model evaluations. Despite this computational overhead for an additional forward pass used to compute multi-token policy ratios after obtaining full trajectories, MPO delivers consistent accuracy improvements on both GSM8K and MATH benchmarks over PPO and GRPO. These results indicate that the added MTP modules contribute minimal training burden while offering substantial efficiency and performance benefits.

Design of Warmup Stage. To ensure stable optimization, we include a brief warm-up phase to pre-train the MTP modules (approximately 12 minutes for one epoch) before the MPO fine-tuning stage (around 7 hours for 20 epochs), incurring only a lightweight additional time cost. This design choice follows prior multi-head decoding frameworks (Cai et al., 2024; Ankner et al., 2024) and represents a standard stabilization step rather than additional complexity unique to our method.

MPO emphasizes a modular and practical architecture; this design enables the framework to improve optimization stability and reasoning performance; furthermore, future work may explore even lighter-weight MTP variants for large-scale deployment (also discussed in Limitations).

992 **E Usage of Large Language Models**

993 During the writing of this paper, large language
994 models were used solely as tools for:

- 995 • **Grammar checking:** Asking the LLM to
996 check whether there is any grammar mistakes
997 in the given paragraph.
- 998 • **Translation:** Asking the LLM to provide
999 a proper translation for certain phrases and
1000 expressions.
- 1001 • **Readability improvement:** Asking the LLM
1002 to give suggestions of how to change the tone
1003 of our written paragraphs for fluency.

1004 **F Reproducibility**

1005 We will briefly go through the tools and data we
1006 use in this paper for reproducibility.

1007 **MTP Modules.** We developed our Multi-token
1008 Prediction modules based on transformers (ver
1009 4.51.1) and trl (ver 0.22.0), both frameworks are
1010 open-sourced by the HuggingFace community. We
1011 modify the modeling library file of the transformers
1012 framework to achieve multi-token prediction.

1013 **MPO Approach.** The overall MPO approach is
1014 developed based on the open-source reinforcement
1015 learning framework Verl (volcano engine, ver
1016 0.3.1). We modify the PPO Trainer Class and Actor
1017 workers Class library to compute and gather multi-
1018 token log probability distributions when computing
1019 the policy loss of MPO.

1020 **Datasets and models.** We conduct our
1021 experiments with open-source datasets GSM8K
1022 and MATH for MTP warmup and MPO training.
1023 The models are available on the open-source AI
1024 community websites.

1025 **Code and implementations.** In the supplemen-
1026 tary materials, we include a detailed README
1027 file that guides the reproduction process, and we
1028 provide our implementation of the causal MTP
1029 module as well as the core code computing MPO
1030 policy loss for reference.