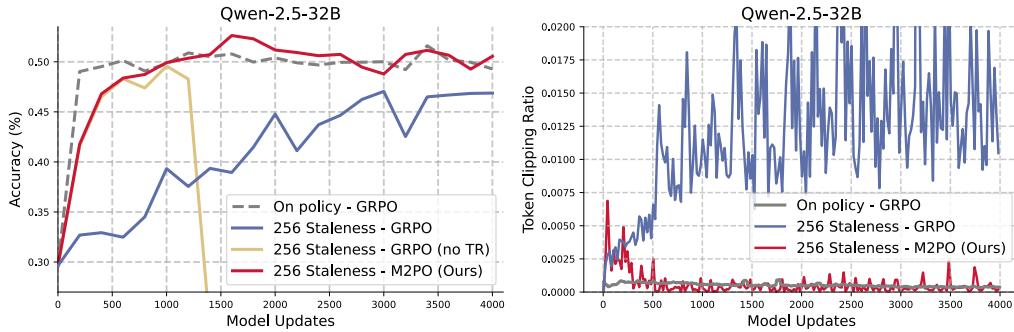


000 001 002 003 004 005 PROSPERITY BEFORE COLLAPSE: HOW FAR CAN OFF- 006 POLICY RL REACH WITH STALE DATA ON LLMs? 007 008 009

010 **Anonymous authors**
011 Paper under double-blind review
012
013
014
015
016
017
018
019
020
021
022
023
024
025

ABSTRACT

011 Reinforcement learning has been central to recent advances in large language
012 model reasoning, but most algorithms rely on on-policy training that demands
013 fresh rollouts at every update, limiting efficiency and scalability. Asynchronous
014 RL systems alleviate this by decoupling rollout generation from training, yet their
015 effectiveness hinges on tolerating large staleness in rollout data, a setting where
016 existing methods either degrade in performance or collapse. We revisit this chal-
017 lenge and uncover a *prosperity-before-collapse* phenomenon: stale data can be as
018 informative as on-policy data if exploited properly. Building on this insight, we
019 introduce **M2PO** (Second-Moment Trust Policy Optimization), which constrains
020 the second moment of importance weights to suppress only extreme outliers while
021 preserving informative updates. Notably, M2PO sharply reduces the fraction of
022 clipped tokens under high staleness (from 1.22% to 0.06% over training), pre-
023 cisely masking high-variance tokens while maintaining stable optimization. Ex-
024 tensive evaluation across six model scales (from 1.7B to 32B) and eight reasoning
025 benchmarks shows that M2PO delivers stable off-policy training even with data
026 stale by at least 256 model updates and matches on-policy performance.
027



038
039 Figure 1: Comparison of on-policy GRPO and off-policy training under a staleness of 256 model
040 updates on Qwen-2.5-32B. **Left:** Standard GRPO suffers from degradation with stale rollouts, while
041 removing the trust region (GRPO no TR) reveals a clear *prosperity-before-collapse* phenomenon.
042 In contrast, M2PO achieves stable training and matches on-policy performance even under high
043 staleness. **Right:** Token clipping ratio comparison shows that M2PO dramatically reduces clipping
044 events compared to GRPO with the same staleness, while avoiding training collapse.

1 INTRODUCTION

047 Reinforcement learning (RL) has been central to recent advances in large language model (LLM)
048 reasoning, driving breakthroughs in systems like OpenAI’s o1 (OpenAI et al., 2024) and DeepSeek’s
049 R1 (DeepSeek-AI et al., 2025; Team et al., 2025). Most existing RL algorithms (Schulman et al.,
050 2017; Zheng et al., 2025b; Yu et al., 2025) for LLMs adopt an on-policy design, as it provides stable
051 training and reliable performance, but the strict requirement for fresh (or limited-staleness) rollouts
052 at every update. **This constraint significantly limits scalability and becomes increasingly imprac-**
053 **tical as we move toward harder, more complex reasoning or agentic tasks, where rollouts involve**
multi-step tool use, expensive program execution, and long reward computation. In such settings,

054 rollout latency can easily range from minutes to hours, making substantial staleness not only com-
 055 mon but often unavoidable. For example, SWE-bench Jimenez et al. (2023) with OpenHands Wang
 056 et al. (2024), rollout latency can become extremely large. In a single run, even with a small batch
 057 size, the end-to-end inference time (including tool call and code execution) can exceed 100 min-
 058 utes, with more than 80 iterations in the environment, which makes on-policy training extremely
 059 inefficient. To overcome this bottleneck, a growing line of RL systems (Fu et al., 2025; Zhu et al.,
 060 2025; Noukhovitch et al., 2024; Zhong et al., 2025; He et al., 2025) have explored asynchronous
 061 designs that decouple rollout from training. Such approaches improve resource utilization and en-
 062 able training to scale more efficiently across large and heterogeneous clusters, but their effectiveness
 063 fundamentally relies on the ability of RL algorithms to tolerate rollout staleness without sacrificing
 064 stability or performance.

065 However, under large rollout staleness, existing RL algorithms struggle to strike the right balance.
 066 Some methods (Schulman et al., 2017; Shao et al., 2024; Zheng et al., 2025a) can maintain stability,
 067 but they often suffer from noticeable performance degradation. Conversely, approaches designed to
 068 maximize performance (Fu et al., 2025; Chen et al., 2025; Su et al., 2025) tend to compromise sta-
 069 bility, frequently leading to training collapse. On-policy methods provide both stability and strong
 070 performance, but their reliance on fresh or only slightly stale rollouts at every update imposes rigid
 071 constraints that hinder scalability. Consequently, an ideal off-policy RL algorithm for LLMs should
 072 enable effective reuse of trajectories collected under outdated policies to preserve strong per-
 073 formance under significant staleness, and ensure stable training that converges competitively with on-
 074 policy methods. Meeting these requirements is key to realizing off-policy RL as a truly scalable
 075 solution for aligning and fine-tuning large language models.

076 In this paper, we aim to investigate the underlying reasons for the limitations of off-policy RL in
 077 LLMs and to design an effective algorithm that fully leverages stale data to unlock its potential.
 078 We begin by revealing an intriguing *Prosperity before Collapse* phenomenon (Yellow curve in Fig-
 079 ure 1 (left)): although RL training without a trust region eventually collapses on stale data, it initially
 080 achieves substantially higher performance than vanilla GRPO with ϵ -clipping. In some cases, it even
 081 matches the performance of the on-policy baseline. From this, we draw an important observation:
 082 *stale data can be as informative as data collected on-policy in RL for LLMs*, but the key chal-
 083 lenge lies in how existing algorithms exploit it. In particular, vanilla GRPO performs poorly under
 084 staleness because stale-data training exhibits a substantially higher clipping rate, with many of the
 085 clipped updates occurring on informative high-entropy tokens (see Figure 4). This disproportionate
 086 clipping on crucial tokens hinders the full utilization of stale training data.

087 This pivotal token masking observation reveals that these high-entropy tokens play a dual role: they
 088 provide the most informative training signal but also introduce the greatest instability under sta-
 089 leness. Therefore, the key challenge is to retain as much learning signal from these tokens as possible
 090 without risking training collapse. Motivated by this, we propose **M2PO** (Second-Moment Trust
 091 Policy Optimization), a novel off-policy RL algorithm that constrains the second moment of impor-
 092 tance weights. Unlike standard ϵ -clipping, which disproportionately suppresses high-entropy tokens
 093 and discards valuable learning signals, M2PO leverages the second-moment metric M_2 . This metric
 094 is both variance-sensitive, capturing instability introduced by high-entropy tokens, and statistically
 095 stable, avoiding the cancellation issues inherent to KL-based measures. By regularizing training at
 096 the batch level through M_2 , M2PO masks only extreme outliers while preserving the majority of
 097 informative updates. As a result, M2PO enables stable off-policy reinforcement learning with stale
 098 data, matching on-policy performance even under large staleness.

099 As illustrated in Figure 1 (left), even when trained exclusively on data stale by at least 256 model
 100 updates, M2PO achieves accuracy comparable to the on-policy baseline (red curve), demon-
 101 strating its ability to fully exploit stale data without sacrificing stability. M2PO achieves this through a more
 102 accurate and adaptive clipping strategy that clips substantially fewer tokens while maintaining train-
 103 ing stability. As shown in Figure 1 (right), M2PO dramatically reduces the fraction of clipped tokens
 104 under high staleness (from 1.22% to 0.06% over the entire training process, see Figure 7b), thereby
 105 preserving more useful training information in stale data. To further validate M2PO effectiveness,
 106 we conduct an extensive evaluation of M2PO across six model scales (ranging from 1.7B to 32B)
 107 and eight math reasoning benchmarks in Section 6. The results show that M2PO consistently deliv-
 108 ers strong performance across all training settings. M2PO also shows insensitivity to the choice of
 109 threshold, with a single value across all experiments, demonstrating its practicality and robustness.

108

2 RELATED WORK

109

110

RLVR. Recent advances (DeepSeek-AI et al., 2025; Yu et al., 2025; Team et al., 2025; Gao et al.,
111 2024) in LLM reasoning show that Reinforcement Learning with Verifiable Reward (RLVR), which
112 relies on verifiable reward signals instead of model-generated scores, can effectively improve model
113 reasoning ability. These gains are achieved using various policy optimization methods such as
114 PPO (Ouyang et al., 2022) and GRPO (Shao et al., 2024). Encouraged by the success of RLVR,
115 a growing body of work (Kazemnejad et al., 2024; Yuan et al., 2025b;a; Yu et al., 2025; Liu et al.,
116 2025b; Luo et al., 2025a; Zhang et al., 2025; Hu, 2025; Xiong et al., 2025) has emerged to fur-
117 ther improve reinforcement learning methods for LLM reasoning. For instance, methods such as
118 VinePPO (Kazemnejad et al., 2024), VC-PPO (Yuan et al., 2025b), and VAPO (Yuan et al., 2025a)
119 aim to enhance LLM reasoning by optimizing the value function.

120

Trust Region in RLVR. While RLVR has been widely adopted for fine-tuning LLMs, a key chal-
121 lenge lies in how to effectively constrain the trust region, not only to stabilize training but also to
122 achieve better learning efficiency and overall performance. To address this, a growing line of work
123 has proposed various strategies to control the policy update, ranging from ratio clipping (Yu et al.,
124 2025), approximate trust region (Fu et al., 2025), sequence-level clipping (Zheng et al., 2025a),
125 asymmetric trust region (Roux et al., 2025; Arnal et al., 2025), and gradient-preserving clipping (Su
126 et al., 2025; Chen et al., 2025). For instance, AREAL (Fu et al., 2025) uses a more recent approxi-
127 mate policy to decide the trust region rather than the behavior model. GSPO (Zheng et al., 2025a)
128 moves from token-level to sequence-level clipping by defining importance ratios on sequence like-
129 lihood. While these methods improve RLVR under moderate settings, most of them focus on rel-
130 atively limited intra-iteration staleness (e.g., 8 or 16) and have not been thoroughly studied under
131 larger off-policy gaps, like extreme staleness. In this work, our goal is to better understand the role
132 of staleness in RLVR and to seek more effective ways of constraining the trust region in RLVR.

133

3 BACKGROUND

134

135

3.1 GROUP RELATIVE POLICY OPTIMIZATION (GRPO)

136

137

Group Relative Policy Optimization (GRPO) (Shao et al., 2024) is a variant of Proximal Policy Opti-
138 mization (PPO) (Ouyang et al., 2022) tailored for language model fine-tuning. Instead of computing
139 advantages using a value function, GRPO normalizes reward scores within groups of responses sam-
140 pled for the same prompt, which largely improves the training efficiency, and aims to maximize the
141 following objective:

142
$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$
143
$$\frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{behav}}(o_i|q)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{behav}}(o_i|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) \right), \quad (1)$$
144

145 where A_i is the advantage, computed using a group of rewards $\{r_1, r_2, \dots, r_G\}$ corresponding to
146 the outputs within each group:

147
$$A_{i,t} = \frac{r_i - \text{mean}(\{R_i\}_{i=1}^G)}{\text{std}(\{R_i\}_{i=1}^G)}. \quad (2)$$
148

149 Similar to PPO, GRPO employs a clipping mechanism to stabilize updates. The ratio $r_i = \frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}$
150 is clipped to $[1 - \epsilon, 1 + \epsilon]$, so that when $A_i > 0$ the policy cannot increase probability mass exces-
151 sively, and when $A_i < 0$ it cannot over-penalize. This prevents large, unstable updates while still
152 allowing normalized group advantages to guide learning.

153

3.2 PERFORMANCE DEGRADATION FROM TRAINING WITH STALE DATA

154

155 To investigate the impact of stale data on reinforcement learning for large language models, we
156 introduce Stale- k RL training, where the model is trained using data generated k model updates

earlier in each training iteration. More specifically, in our training setup, each training step consists of four model updates, a configuration commonly used in recent work (Zheng et al., 2025b; Wang et al., 2025b; Yu et al., 2025; Chen et al., 2025; Zheng et al., 2025a). Thus, even stale-0 ($s=0$) training has a staleness between 0 and 3. stale-256 ($s=256$) training has a staleness between 256 and 259). During the first k model updates, since no stale model is yet available, the model is trained on data generated by the original base model, with different training data used in each iteration. In this setup, all training data after the initial phase comes from stale models, allowing us to study how stale data affects the dynamics and effectiveness of RL training.

Stale- k RL training. In our “stale- k ” setup, the training batch at step t is collected k updates earlier, but we implement this by storing the generated data rather than storing old checkpoints. Concretely, at update step $t - k$ we use the current policy to generate rollouts and place them into a buffer, and after k further updates, these trajectories are consumed once for training at step t . Each sampled trajectory is used exactly once and is never reused across multiple updates. More training details can be found at Appendix B.

As shown in Figure 2, we train Qwen2.5-Math-7B (Yang et al., 2024) with GRPO under varying staleness levels and report test accuracy. The results reveal a clear trend: as staleness increases, model performance degrades and convergence slows. In particular, low-staleness training achieves higher accuracy, whereas high-staleness training converges more slowly to lower performance.

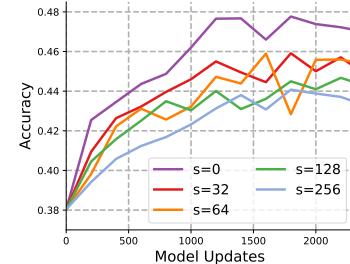


Figure 2: RL with stale data on Qwen2.5-Math-7B, reporting avg. accuracy on all benchmarks.

4 PROSPERITY BEFORE COLLAPSE: STALE DATA CONTAIN ENOUGH TRAINING INFORMATION IN RL ON LLMs

In this section, we investigate why RL on LLM deteriorates when trained on stale data generated by earlier policies.

First, we reveal an intriguing *prosperity-before-collapse* phenomenon: although off-policy RL training without a trust region eventually collapses on stale data, it achieves substantially higher performance than GRPO with ϵ -clipping before collapse, even matching on-policy results. Next, we study the causes of GRPO’s inferior performance when trained with stale data.

Prosperity before collapse: training without a trust region. To disentangle whether the performance drop stems from stale data generated by highly shifted old policies or from biases introduced by the training algorithm, we remove the trust region entirely to remove bias from the training algorithm. Surprisingly, we observe a distinct *prosperity-before-collapse* phenomenon. As shown in Figure 1 and Figure 3, although training without a trust region eventually collapses, it achieves substantially better performance prior to collapse. In fact, under stale data ($s=256$), the no-clipping setting initially outperforms clipped training, sometimes even matching on-policy baselines.

Pivotal token masking by ϵ -clipping when training with stale data. As also discussed in recent work (Su et al., 2025; Chen et al., 2025), ϵ -clipping may inadvertently mask important tokens, preventing them from contributing useful training signals. We extend this observation to the asynchronous setting and show that the problem becomes substantially more severe when training with stale data, since larger staleness induces a greater mismatch between the behavior and target poli-

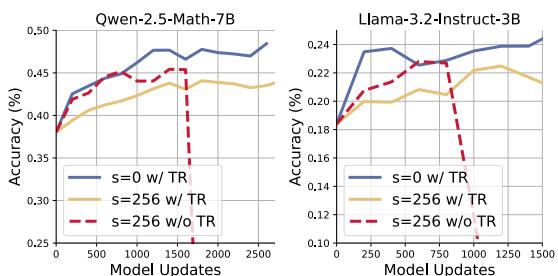


Figure 3: Prosperity before Collapse. Training without a trust region (TR) ($\epsilon = \infty$) under stale data ($s = 256$) initially achieves higher performance than clipped training, sometimes even matching the on-policy baseline ($s = 0$). However, it eventually collapses due to uncontrolled variance.

Figure 3 shows two plots for 'Owen-2.5-Math-7B' and 'Llama-3.2-Instruct-3B'. Both plots show accuracy (%) on the y-axis (0.25 to 0.50) versus model updates on the x-axis. Three lines are shown in each: s=0 w/ TR (solid blue), s=256 w/ TR (solid orange), and s=256 w/o TR (dashed red). In both cases, the s=0 w/ TR and s=256 w/ TR lines start at ~0.40 accuracy and rise to ~0.45 accuracy by 1000 updates. The s=256 w/o TR lines start at ~0.40 accuracy and rise more slowly, reaching ~0.43 accuracy by 1000 updates. After 1000 updates, the s=256 w/o TR lines drop sharply to ~0.35 accuracy by 1500 updates, indicating collapse. The right graph shows a slight increase in accuracy for the s=256 w/o TR line after 1500 updates, reaching ~0.38 accuracy by 2500 updates.

216 cies. As illustrated in Figure 4a, the clipping ratio increases sharply under large staleness ($s = 256$),
 217 while remaining negligible in the on-policy baseline.
 218

219 To better understand this phenomenon,
 220 we conduct a quantitative analysis on 90
 221 million training tokens collected during
 222 Qwen2.5-Math-7B training with staleness
 223 256. Specifically, we gather all training
 224 tokens generated between 800 and 1200
 225 model updates, ensuring the model is
 226 already in a stable training phase but
 227 before convergence. Figure 4b shows a clear
 228 trend: as $|r - 1|$ increases, the average to-
 229 ken entropy also rises.

230 This indicates that ϵ -clipping dis-
 231 proportionately prunes high-entropy tokens,
 232 which are typically the most informa-
 233 tive for model improvement (Wang et al.,
 234 2025a; Gao et al., 2025; Cui et al., 2025).

235 Consequently, clipping under stale data leads to degraded performance.

236 This observation reveals a dilemma: while high-entropy tokens are crucial for learning progress,
 237 they also introduce instability in the off-policy setting, which motivates our key research question:

238 *Can a more accurate and adaptive trust region strategy preserve the benefits of
 239 stale data while ensuring stable training?*

242 5 SECOND-MOMENT TRUST POLICY OPTIMIZATION

244 In this section, we propose Second-Moment Trust Policy Optimization (M2PO), a novel policy
 245 optimization algorithm providing a more effective trust region for off-policy training with stale data.
 246 As discussed in Section 4, as high-entropy tokens are a double-edged sword, the key challenge in
 247 designing an effective trust region algorithm is how to best harness the rich information in high-
 248 entropy tokens without letting them destabilize training.

250 5.1 MEASURING DISTRIBUTION GAP WITH THE SECOND MOMENT

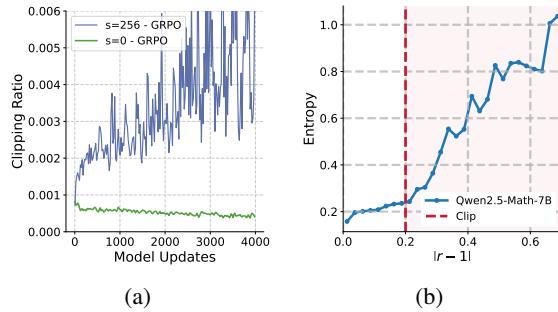
252 The main source of instability in off-policy RL lies in the distributional mismatch between the be-
 253 havior policy that generates training data and the current policy being optimized (Schulman et al.,
 254 2015; 2017). As the divergence between these two distributions grows, importance sampling cor-
 255 rections produce high-variance gradient estimates, leading to noisy and unreliable updates. Our
 256 motivation is therefore to constrain the distributional gap between π_{behav} and π_θ at the batch level,
 257 directly coupling the constraint with model updates while preventing over-constraining of token-
 258 level variations.

259 A natural choice to measure distribution is the batch-level KL divergence, a metric widely adopted
 260 to monitor stability in RL:

$$262 \hat{KL} = \frac{1}{N} \sum_{i=1}^N \hat{KL}_i = -\frac{1}{N} \sum_{i=1}^N \log r_i = -\frac{1}{N} \sum_{i=1}^N \log \frac{\pi_\theta(a_i | s_i)}{\pi_{\text{behav}}(a_i | s_i)}, \quad (3)$$

264 where N is the number of tokens in a batch.

266 However, batch-level KL suffers from two key limitations. First, because it is computed from single-
 267 sample estimates, individual \hat{KL}_i can be positive or negative, leading to *cancellation effects* where
 268 large deviations offset each other and produce deceptively small KL values. Second, tokens with
 269 large ratios ($r_i > 1$) are not properly constrained, as their negative \hat{KL}_i actually decreases the
 estimated KL, even though such tokens can contribute to training instability (Schulman et al., 2017).



238 Figure 4: (a) Clipping ratio dynamics during RL train-
 239 ing on the Qwen-2.5-Math-7B model. (b) Relation-
 240 ship between average token entropy and the distance
 241 between the importance sampling ratio r and 1.

270 To overcome these limitations, we propose to use the second moment of the log-ratio to measure the
 271 distribution gap between behavior and current policy. Formally, we define
 272

$$273 \quad 274 \quad \hat{M}_2 = \frac{1}{N} \sum_{i=1}^N \hat{M}_{2,i} = \frac{1}{N} \sum_{i=1}^N (\log r_i)^2 = \frac{1}{N} \sum_{i=1}^N [\log \frac{\pi_\theta(a_i | s_i)}{\pi_{\text{behav}}(a_i | s_i)}]^2, \quad 275 \quad 276 \quad 277 \quad 278 \quad 279 \quad 280 \quad 281 \quad 282 \quad 283 \quad 284 \quad 285 \quad 286 \quad 287 \quad 288 \quad 289 \quad 290 \quad 291 \quad 292 \quad 293 \quad 294 \quad 295 \quad 296 \quad 297 \quad 298 \quad 299 \quad 300 \quad 301 \quad 302 \quad 303 \quad 304 \quad 305 \quad 306 \quad 307 \quad 308 \quad 309 \quad 310 \quad 311 \quad 312 \quad 313 \quad 314 \quad 315 \quad 316 \quad 317 \quad 318 \quad 319 \quad 320 \quad 321 \quad 322 \quad 323$$

This choice is motivated by two key advantages of \hat{M}_2 over the batch \hat{KL} . First, each per-token estimate $\hat{M}_{2,i} = (\log r_i)^2$ is always non-negative, so the constraint can be reliably applied even when $r > 1$. Second, while the batch KL only measures the mean shift between policies, M_2 also reflects the variance of importance weights. This makes \hat{M}_2 more sensitive to outliers and noisy tokens with extreme ratios r_i^1 .

Furthermore, Theorem 1 shows that although M_2 does not directly constrain $r - 1$ like ϵ clipping, it nevertheless provides an upper bound on the Pearson chi-square divergence $\mathbb{E}[(r - 1)^2]$ between the new and behavior policies. The proof is provided in Appendix D.

Theorem 1 (Bounding χ^2 by M_2). *Let $r = \frac{\pi_{\text{new}}}{\pi_{\text{behav}}}$ be the importance ratio and assume $1/R \leq r \leq R$. Define the log-ratio second moment*

$$M_2 = \mathbb{E}_{a \sim \pi_{\text{behav}}}[(\log r(a))^2].$$

Let the Pearson chi-square divergence between π_{new} and π_{behav} be

$$\chi^2(\pi_{\text{new}} \| \pi_{\text{behav}}) = \mathbb{E}_{a \sim \pi_{\text{behav}}} \left[\left(\frac{\pi_{\text{new}}(a)}{\pi_{\text{behav}}(a)} - 1 \right)^2 \right] = \mathbb{E}_{\pi_{\text{behav}}}[(r - 1)^2].$$

Then

$$\chi^2(\pi_{\text{new}} \| \pi_{\text{behav}}) \leq R^2 M_2.$$

5.2 SECOND-MOMENT TRUST POLICY OPTIMIZATION

As illustrated in Algorithm 1, to maintain training stability, M2PO applies a masking strategy that selectively excludes tokens until the batch-level \hat{M}_2 of the remaining tokens falls below a predefined threshold τ_{M_2} . Importantly, we observe that τ_{M_2} is not a sensitive hyperparameter (see Figure 8). Across all our experiments, we consistently set $\tau_{M_2} = 0.04$, and this single setting proved effective for stabilizing training in all training scenarios.

Only constrain trust-region tokens. Although the PPO loss clips the ratio on both the upper and lower sides, due to the use of the \min operator, not all tokens are actually clipped. In practice, clipping only occurs for tokens where $A > 0$ and $r > 1$, or $A < 0$ and $r < 1$. Following the PPO setting, we therefore apply the M_2 constraint exclusively to tokens that satisfy these conditions. Finally, with the result mask M , we update the policy by maximizing the following objective²:

$$\mathcal{J}_{\text{M2PO}}(\theta) = \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} M_{i,t} \frac{\pi_\theta(o_i | q)}{\pi_{\theta_{\text{behav}}}(o_i | q)} A_{i,t}, \quad M_{i,t} \in \{0, 1\}, \quad 315 \quad 316 \quad 317 \quad 318 \quad 319 \quad 320 \quad 321 \quad 322 \quad 323 \quad 324 \quad 325 \quad 326 \quad 327 \quad 328 \quad 329 \quad 330 \quad 331 \quad 332 \quad 333 \quad 334 \quad 335 \quad 336 \quad 337 \quad 338 \quad 339 \quad 340 \quad 341 \quad 342 \quad 343 \quad 344 \quad 345 \quad 346 \quad 347 \quad 348 \quad 349 \quad 350 \quad 351 \quad 352 \quad 353 \quad 354 \quad 355 \quad 356 \quad 357 \quad 358 \quad 359 \quad 360 \quad 361 \quad 362 \quad 363 \quad 364 \quad 365 \quad 366 \quad 367 \quad 368 \quad 369 \quad 370 \quad 371 \quad 372 \quad 373 \quad 374 \quad 375 \quad 376 \quad 377 \quad 378 \quad 379 \quad 380 \quad 381 \quad 382 \quad 383 \quad 384 \quad 385 \quad 386 \quad 387 \quad 388 \quad 389 \quad 390 \quad 391 \quad 392 \quad 393 \quad 394 \quad 395 \quad 396 \quad 397 \quad 398 \quad 399 \quad 400 \quad 401 \quad 402 \quad 403 \quad 404 \quad 405 \quad 406 \quad 407 \quad 408 \quad 409 \quad 410 \quad 411 \quad 412 \quad 413 \quad 414 \quad 415 \quad 416 \quad 417 \quad 418 \quad 419 \quad 420 \quad 421 \quad 422 \quad 423 \quad 424 \quad 425 \quad 426 \quad 427 \quad 428 \quad 429 \quad 430 \quad 431 \quad 432 \quad 433 \quad 434 \quad 435 \quad 436 \quad 437 \quad 438 \quad 439 \quad 440 \quad 441 \quad 442 \quad 443 \quad 444 \quad 445 \quad 446 \quad 447 \quad 448 \quad 449 \quad 450 \quad 451 \quad 452 \quad 453 \quad 454 \quad 455 \quad 456 \quad 457 \quad 458 \quad 459 \quad 460 \quad 461 \quad 462 \quad 463 \quad 464 \quad 465 \quad 466 \quad 467 \quad 468 \quad 469 \quad 470 \quad 471 \quad 472 \quad 473 \quad 474 \quad 475 \quad 476 \quad 477 \quad 478 \quad 479 \quad 480 \quad 481 \quad 482 \quad 483 \quad 484 \quad 485 \quad 486 \quad 487 \quad 488 \quad 489 \quad 490 \quad 491 \quad 492 \quad 493 \quad 494 \quad 495 \quad 496 \quad 497 \quad 498 \quad 499 \quad 500 \quad 501 \quad 502 \quad 503 \quad 504 \quad 505 \quad 506 \quad 507 \quad 508 \quad 509 \quad 510 \quad 511 \quad 512 \quad 513 \quad 514 \quad 515 \quad 516 \quad 517 \quad 518 \quad 519 \quad 520 \quad 521 \quad 522 \quad 523 \quad 524 \quad 525 \quad 526 \quad 527 \quad 528 \quad 529 \quad 530 \quad 531 \quad 532 \quad 533 \quad 534 \quad 535 \quad 536 \quad 537 \quad 538 \quad 539 \quad 540 \quad 541 \quad 542 \quad 543 \quad 544 \quad 545 \quad 546 \quad 547 \quad 548 \quad 549 \quad 550 \quad 551 \quad 552 \quad 553 \quad 554 \quad 555 \quad 556 \quad 557 \quad 558 \quad 559 \quad 560 \quad 561 \quad 562 \quad 563 \quad 564 \quad 565 \quad 566 \quad 567 \quad 568 \quad 569 \quad 570 \quad 571 \quad 572 \quad 573 \quad 574 \quad 575 \quad 576 \quad 577 \quad 578 \quad 579 \quad 580 \quad 581 \quad 582 \quad 583 \quad 584 \quad 585 \quad 586 \quad 587 \quad 588 \quad 589 \quad 590 \quad 591 \quad 592 \quad 593 \quad 594 \quad 595 \quad 596 \quad 597 \quad 598 \quad 599 \quad 600 \quad 601 \quad 602 \quad 603 \quad 604 \quad 605 \quad 606 \quad 607 \quad 608 \quad 609 \quad 610 \quad 611 \quad 612 \quad 613 \quad 614 \quad 615 \quad 616 \quad 617 \quad 618 \quad 619 \quad 620 \quad 621 \quad 622 \quad 623 \quad 624 \quad 625 \quad 626 \quad 627 \quad 628 \quad 629 \quad 630 \quad 631 \quad 632 \quad 633 \quad 634 \quad 635 \quad 636 \quad 637 \quad 638 \quad 639 \quad 640 \quad 641 \quad 642 \quad 643 \quad 644 \quad 645 \quad 646 \quad 647 \quad 648 \quad 649 \quad 650 \quad 651 \quad 652 \quad 653 \quad 654 \quad 655 \quad 656 \quad 657 \quad 658 \quad 659 \quad 660 \quad 661 \quad 662 \quad 663 \quad 664 \quad 665 \quad 666 \quad 667 \quad 668 \quad 669 \quad 670 \quad 671 \quad 672 \quad 673 \quad 674 \quad 675 \quad 676 \quad 677 \quad 678 \quad 679 \quad 680 \quad 681 \quad 682 \quad 683 \quad 684 \quad 685 \quad 686 \quad 687 \quad 688 \quad 689 \quad 690 \quad 691 \quad 692 \quad 693 \quad 694 \quad 695 \quad 696 \quad 697 \quad 698 \quad 699 \quad 700 \quad 701 \quad 702 \quad 703 \quad 704 \quad 705 \quad 706 \quad 707 \quad 708 \quad 709 \quad 710 \quad 711 \quad 712 \quad 713 \quad 714 \quad 715 \quad 716 \quad 717 \quad 718 \quad 719 \quad 720 \quad 721 \quad 722 \quad 723 \quad 724 \quad 725 \quad 726 \quad 727 \quad 728 \quad 729 \quad 730 \quad 731 \quad 732 \quad 733 \quad 734 \quad 735 \quad 736 \quad 737 \quad 738 \quad 739 \quad 740 \quad 741 \quad 742 \quad 743 \quad 744 \quad 745 \quad 746 \quad 747 \quad 748 \quad 749 \quad 750 \quad 751 \quad 752 \quad 753 \quad 754 \quad 755 \quad 756 \quad 757 \quad 758 \quad 759 \quad 760 \quad 761 \quad 762 \quad 763 \quad 764 \quad 765 \quad 766 \quad 767 \quad 768 \quad 769 \quad 770 \quad 771 \quad 772 \quad 773 \quad 774 \quad 775 \quad 776 \quad 777 \quad 778 \quad 779 \quad 780 \quad 781 \quad 782 \quad 783 \quad 784 \quad 785 \quad 786 \quad 787 \quad 788 \quad 789 \quad 790 \quad 791 \quad 792 \quad 793 \quad 794 \quad 795 \quad 796 \quad 797 \quad 798 \quad 799 \quad 800 \quad 801 \quad 802 \quad 803 \quad 804 \quad 805 \quad 806 \quad 807 \quad 808 \quad 809 \quad 810 \quad 811 \quad 812 \quad 813 \quad 814 \quad 815 \quad 816 \quad 817 \quad 818 \quad 819 \quad 820 \quad 821 \quad 822 \quad 823 \quad 824 \quad 825 \quad 826 \quad 827 \quad 828 \quad 829 \quad 830 \quad 831 \quad 832 \quad 833 \quad 834 \quad 835 \quad 836 \quad 837 \quad 838 \quad 839 \quad 840 \quad 841 \quad 842 \quad 843 \quad 844 \quad 845 \quad 846 \quad 847 \quad 848 \quad 849 \quad 850 \quad 851 \quad 852 \quad 853 \quad 854 \quad 855 \quad 856 \quad 857 \quad 858 \quad 859 \quad 860 \quad 861 \quad 862 \quad 863 \quad 864 \quad 865 \quad 866 \quad 867 \quad 868 \quad 869 \quad 870 \quad 871 \quad 872 \quad 873 \quad 874 \quad 875 \quad 876 \quad 877 \quad 878 \quad 879 \quad 880 \quad 881 \quad 882 \quad 883 \quad 884 \quad 885 \quad 886 \quad 887 \quad 888 \quad 889 \quad 890 \quad 891 \quad 892 \quad 893 \quad 894 \quad 895 \quad 896 \quad 897 \quad 898 \quad 899 \quad 900 \quad 901 \quad 902 \quad 903 \quad 904 \quad 905 \quad 906 \quad 907 \quad 908 \quad 909 \quad 910 \quad 911 \quad 912 \quad 913 \quad 914 \quad 915 \quad 916 \quad 917 \quad 918 \quad 919 \quad 920 \quad 921 \quad 922 \quad 923 \quad 924 \quad 925 \quad 926 \quad 927 \quad 928 \quad 929 \quad 930 \quad 931 \quad 932 \quad 933 \quad 934 \quad 935 \quad 936 \quad 937 \quad 938 \quad 939 \quad 940 \quad 941 \quad 942 \quad 943 \quad 944 \quad 945 \quad 946 \quad 947 \quad 948 \quad 949 \quad 950 \quad 951 \quad 952 \quad 953 \quad 954 \quad 955 \quad 956 \quad 957 \quad 958 \quad 959 \quad 960 \quad 961 \quad 962 \quad 963 \quad 964 \quad 965 \quad 966 \quad 967 \quad 968 \quad 969 \quad 970 \quad 971 \quad 972 \quad 973 \quad 974 \quad 975 \quad 976 \quad 977 \quad 978 \quad 979 \quad 980 \quad 981 \quad 982 \quad 983 \quad 984 \quad 985 \quad 986 \quad 987 \quad 988 \quad 989 \quad 990 \quad 991 \quad 992 \quad 993 \quad 994 \quad 995 \quad 996 \quad 997 \quad 998 \quad 999 \quad 1000 \quad 1001 \quad 1002 \quad 1003 \quad 1004 \quad 1005 \quad 1006 \quad 1007 \quad 1008 \quad 1009 \quad 1010 \quad 1011 \quad 1012 \quad 1013 \quad 1014 \quad 1015 \quad 1016 \quad 1017 \quad 1018 \quad 1019 \quad 1020 \quad 1021 \quad 1022 \quad 1023 \quad 1024 \quad 1025 \quad 1026 \quad 1027 \quad 1028 \quad 1029 \quad 1030 \quad 1031 \quad 1032 \quad 1033 \quad 1034 \quad 1035 \quad 1036 \quad 1037 \quad 1038 \quad 1039 \quad 1040 \quad 1041 \quad 1042 \quad 1043 \quad 1044 \quad 1045 \quad 1046 \quad 1047 \quad 1048 \quad 1049 \quad 1050 \quad 1051 \quad 1052 \quad 1053 \quad 1054 \quad 1055 \quad 1056 \quad 1057 \quad 1058 \quad 1059 \quad 1060 \quad 1061 \quad 1062 \quad 1063 \quad 1064 \quad 1065 \quad 1066 \quad 1067 \quad 1068 \quad 1069 \quad 1070 \quad 1071 \quad 1072 \quad 1073 \quad 1074 \quad 1075 \quad 1076 \quad 1077 \quad 1078 \quad 1079 \quad 1080 \quad 1081 \quad 1082 \quad 1083 \quad 1084 \quad 1085 \quad 1086 \quad 1087 \quad 1088 \quad 1089 \quad 1090 \quad 1091 \quad 1092 \quad 1093 \quad 1094 \quad 1095 \quad 1096 \quad 1097 \quad 1098 \quad 1099 \quad 1100 \quad 1101 \quad 1102 \quad 1103 \quad 1104 \quad 1105 \quad 1106 \quad 1107 \quad 1108 \quad 1109 \quad 1110 \quad 1111 \quad 1112 \quad 1113 \quad 1114 \quad 1115 \quad 1116 \quad 1117 \quad 1118 \quad 1119 \quad 1120 \quad 1121 \quad 1122 \quad 1123 \quad 1124 \quad 1125 \quad 1126 \quad 1127 \quad 1128 \quad 1129 \quad 1130 \quad 1131 \quad 1132 \quad 1133 \quad 1134 \quad 1135 \quad 1136 \quad 1137 \quad 1138 \quad 1139 \quad 1140 \quad 1141 \quad 1142 \quad 1143 \quad 1144 \quad 1145 \quad 1146 \quad 1147 \quad 1148 \quad 1149 \quad 1150 \quad 1151 \quad 1152 \quad 1153 \quad 1154 \quad 1155 \quad 1156 \quad 1157 \quad 1158 \quad 1159 \quad 1160 \quad 1161 \quad 1162 \quad 1163 \quad 1164 \quad 1165 \quad 1166 \quad 1167 \quad 1168 \quad 1169 \quad 1170 \quad 1171 \quad 1172 \quad 1173 \quad 1174 \quad 1175 \quad 1176 \quad 1177 \quad 1178 \quad 1179 \quad 1180 \quad 1181 \quad 1182 \quad 1183 \quad 1184 \quad 1185 \quad 1186 \quad 1187 \quad 1188 \quad 1189 \quad 1190 \quad 1191 \quad 1192 \quad 1193 \quad 1194 \quad 1195 \quad 1196 \quad 1197 \quad 1198 \quad 1199 \quad 1200 \quad 1201 \quad 1202 \quad 1203 \quad 1204 \quad 1205 \quad 1206 \quad 1207 \quad 1208 \quad 1209 \quad 1210 \quad 1211 \quad 1212 \quad 1213 \quad 1214 \quad 1215 \quad 1216 \quad 1217 \quad 1218 \quad 1219 \quad 1220 \quad 1221 \quad 1222 \quad 1223 \quad 1224 \quad 1225 \quad 1226 \quad 1227 \quad 1228 \quad 1229 \quad 1230 \quad 1231 \quad 1232 \quad 1233 \quad 1234 \quad 1235 \quad 1236 \quad 1237 \quad 1238 \quad 1239 \quad 1240 \quad 1241 \quad 1242 \quad 1243 \quad 1244 \quad 1245 \quad 1246 \quad 1247 \quad 1248 \quad 1249 \quad 1250 \quad 1251 \quad 1252 \quad 1253 \quad 1254 \quad 1255 \quad 1256 \quad 1257 \quad 1258 \quad 1259 \quad 1260 \quad 1261 \quad 1262 \quad 1263 \quad 1264 \quad 1265 \quad 1266 \quad 1267 \quad 1268 \quad 1269 \quad 1270 \quad 1271 \quad 1272 \quad 1273 \quad 1274 \quad 1275 \quad 1276 \quad 1277 \quad 1278 \quad 1279 \quad 1280 \quad 1281 \quad 1282 \quad 1283 \quad 1284 \quad 1285 \quad 1286 \quad 1287 \quad 1288 \quad 1289 \quad 1290 \quad 1291 \quad 1292 \quad 1293 \quad 1294 \quad 1295 \quad 1296 \quad 1297 \quad 1298 \quad 1299 \quad 1300 \quad 1301 \quad 1302 \quad 1303 \quad 1304 \quad 1305 \quad 1306 \quad 1307 \quad 1308 \quad 1309 \quad 1310 \quad 1311 \quad 1312 \quad 1313 \quad 1314 \quad 1315 \quad 1316 \quad 1317 \quad 1318 \quad 1319 \quad 1320 \quad 1321 \quad 1322 \quad 1323 \quad 1324 \quad 1325 \quad 1326 \quad 1327 \quad 1328 \quad 1329 \quad 1330 \quad 1331 \quad 1332 \quad 1333 \quad 1334 \quad 1335 \quad 1336 \quad 1337 \quad 1338 \quad 1339 \quad 1340 \quad 1341 \quad 1342 \quad 1343 \quad 1344 \quad 1345 \quad 1346 \quad 1347 \quad 1348 \quad 1349 \quad 1350 \quad 1351 \quad 1352 \quad 1353 \quad 1354 \quad 1355 \quad 1356 \quad 1357 \quad 1358 \quad 1359 \quad 1360 \quad 1361 \quad 1362 \quad 1363 \quad 1364 \quad 1365 \quad 1366 \quad 1367 \quad 1368 \quad 1369 \quad 1370 \quad 1371 \quad 1372 \quad 1373 \quad 1374 \quad 1375 \quad 1376 \quad 1377 \quad 1378 \quad 1379 \quad 1380 \quad 1381 \quad 1382 \quad 1383 \quad 1384 \quad 1385 \quad 1386 \quad 1387 \quad 1388 \quad 1389 \quad 1390 \quad 1391 \quad 1392 \quad 1393 \quad 1394 \quad 1395 \quad 1396 \quad 1397 \quad 1398 \quad 1399 \quad 1400 \quad 1401 \quad 1402 \quad 1403 \quad 1404 \quad 1405 \quad 1406 \quad 1407 \quad 1408 \quad 1409 \quad 1410 \quad 1411 \quad 1412 \quad 1413 \quad 1414 \quad 1415 \quad 1416 \quad 1417 \quad 1418 \quad 1419 \quad 1420 \quad 1421 \quad 1422 \quad 1423 \quad 1424 \quad 1425 \quad 1426 \quad 1427 \quad 1428 \quad 1429 \quad 1430 \quad 1431 \quad 1432 \quad 1433 \quad 1434 \quad 1435 \quad 1436 \quad 1437 \quad 1438 \quad 1439 \quad 1440 \quad 1441 \quad 1442 \quad 1443 \quad 1444 \quad 1445 \quad 1446 \quad 1447 \quad 1448 \quad 1449 \quad 1450 \quad 1451 \quad 1452 \quad 1453 \quad 1454 \quad 1455 \quad 1456 \quad 1457 \quad 1458 \quad 1459 \quad 1460 \quad 1461 \quad 1462 \quad 1463 \quad 1464 \quad 1465 \quad 1466 \quad 1467 \quad 1468 \quad 1469 \quad 1470 \quad 1471 \quad 1472 \quad 1473 \quad 1474 \quad 1475 \quad 1476 \quad 1477 \quad 1478 \quad 1479 \quad 1480 \quad 1481 \quad 1482 \quad 1483 \quad 1484 \quad 1485 \quad 1486 \quad 1487 \quad 1488 \quad 1489 \quad 1490 \quad 1491 \quad 1492 \quad 1493 \quad 1494 \quad 1495 \quad 1496 \quad 1497 \quad 1498 \quad 1499 \quad 1500 \quad 1501 \quad 1502 \quad 1503 \quad 1504 \quad 1505 \quad 1506 \quad 1507 \quad 1508 \quad 1509 \quad 1510 \quad 1511 \quad 1512 \quad 1513 \quad 1514 \quad 1515 \quad 1516 \quad 1517 \quad 1518 \quad 1519 \quad 1520 \quad 1521 \quad 1522 \quad 1523 \quad 1524 \quad 1525 \quad 1526 \quad 1527 \quad 1528 \quad 1529 \quad 1530 \quad 1531 \quad 1532 \quad 1533 \quad 1534 \quad 1535 \quad 1536 \quad 1537 \quad 1538 \quad 1539 \quad 1540 \quad 1541 \quad 1542 \quad 1543 \quad 1544 \quad 1545 \quad 1546 \quad 1547 \quad 1548 \quad 1549 \quad 1550 \quad 1551 \quad 1552 \quad 1553 \quad 1554 \quad 1555 \quad 1556 \quad 1557 \quad 1558 \quad 1559 \quad 1560 \quad 1561 \quad 1562 \quad 1563 \quad 1564 \quad 1565 \quad 1566 \quad 1567 \quad 1568 \quad 1569 \quad 1570 \quad 1571 \quad 1572 \quad 1573 \quad 1574 \quad 1575 \quad 1576 \quad 1577 \quad 1578 \quad 1579 \quad 1580 \quad 1581 \quad 1582 \quad 1583 \quad 1584 \quad 1585 \quad 1586 \quad 1587 \quad 1588 \quad 1589 \quad 1590 \quad 1591 \quad 1592 \quad 1593 \quad 1594 \quad 1595 \quad 1596 \quad 1597 \quad 1598 \quad 1599 \quad 1600 \quad 1601 \quad 1602 \quad 1603 \quad 1604 \quad 1605 \quad 1606 \quad 1607 \quad 1608 \quad 1609 \quad 1610 \quad 1611 \quad 1612 \quad 1613 \quad 1614 \quad 1615 \quad 1616 \quad 1617 \quad 1618 \quad 1619 \quad 1620 \quad 1621 \quad 1622 \quad 1623 \quad 1624 \quad 1625 \quad 1626 \quad 1627 \quad 1628 \quad 1629 \quad 1630 \quad 1631 \quad 1632 \quad 1633 \quad 1634 \quad 1635 \quad 1636 \quad 1637 \quad 1638 \quad 1639 \quad 1640 \quad 1641 \quad 1642 \quad 1643 \quad 1644 \quad 1645 \quad 1646 \quad 1647 \quad 1648 \quad 1649 \quad 1650 \quad 1651 \quad 1652 \quad 1653 \quad 1654 \quad 1655 \quad 1656 \quad 1657 \quad 1658 \quad 1659 \quad 1660 \quad 1661 \quad 1662 \quad 1663 \quad 1664 \quad 1665 \quad 1666 \quad 1667 \quad 1668 \quad 1669 \quad 1670 \quad 1671 \quad 1672 \quad 1673 \quad 1674 \quad 1675 \quad 1676 \quad 1677 \quad 1678 \quad 1679 \quad 1680 \quad 1681 \quad 1682 \quad 1683 \quad 1684 \quad 1685 \quad 1686 \quad 1687 \quad 1688 \quad 1689 \quad 1690 \quad 1691 \quad 1692 \quad 1693 \quad 1694 \quad 1695 \quad 1696 \quad 1697 \quad 1698 \quad 1699 \quad 1700 \quad 1701 \quad 1702 \quad 1703 \quad 1704 \quad 1705 \quad 1706 \quad 1707 \quad 1708 \quad 1709 \quad 1710 \quad 1711 \quad 1712 \quad 1713 \quad 1714 \quad 1715 \quad 1716 \quad 1717 \quad 1718 \quad 1719 \quad 1720 \quad 1721 \quad 1722 \quad 1723 \quad 1724 \quad 1725 \quad 1726 \quad 1727 \quad 1728 \quad 1729 \quad 1730 \quad 1731 \quad 1732 \quad 1733 \quad 1734 \quad 1735 \quad 1736 \quad 1737 \quad 1738 \quad 1739 \quad 1740 \quad 1741 \quad 1742 \quad 1743 \quad 1744 \quad 1745 \quad 1746 \quad 1747 \quad 1748 \quad 1749 \quad 1750 \quad 1751 \quad 1752 \quad 1753 \quad 1754 \quad 1755 \quad 1756 \quad 1757 \quad 1758 \quad 1759 \quad 1760 \quad 1761 \quad 176$$

324 **6 EXPERIMENTS**
 325

326 In this section, we present an extensive evaluation across six models (from 1.7B to 32B) on eight
 327 benchmarks. The results demonstrate that, even when trained with extremely stale data, M2PO
 328 achieves performance comparable to on-policy GRPO and significantly outperforms other baselines:
 329

- 330 • In Section 6.2, we show that M2PO achieves **accuracy on par with on-policy baselines** under
 331 large staleness ($s = 256$), and outperforms baselines by up to 11.2% in average accuracy.
- 332 • In Section 6.3, we provide a detailed analysis of how M2PO boosts off-policy RL performance
 333 while preserving training stability. We also show that its sole threshold hyperparameter τ_{M_2}
 334 is insensitive to variation, ensuring ease of use in practice.

335 **6.1 EXPERIMENTAL SETTINGS**
 336

337 **Models & Datasets.** To verify the effectiveness of our method, we extensively evaluate M2PO
 338 on six models: Qwen2.5-Math-7B (Yang et al., 2024), Llama-3.2-3B-Instruct (Dubey et al., 2024),
 339 Qwen3-Base-1.7B/4B/8B (Yang et al., 2025), and Qwen2.5-32B (Yang et al., 2024). For Qwen2.5-
 340 Math-7B, we use a context length of 4k, which is the maximum for this series, while for all other
 341 models the context length is set to 16k. For training, we adopt the DeepScaleR (Luo et al., 2025b)
 342 math dataset.

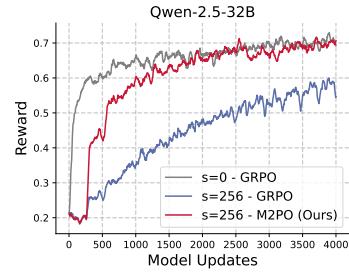
343 **Training & Evaluation.** Our method is implemented based on verl (Sheng et al., 2024) pipeline and
 344 uses vLLM (Kwon et al., 2023) for rollout. We use a mix of H100 and H200 servers for training, de-
 345 pending on resource availability. For benchmark datasets, we use eight widely used complex mathe-
 346 matical reasoning benchmarks to evaluate the performance of trained models: Math500 (Hendrycks
 347 et al., 2021; Lightman et al., 2023), AIME24/25 (Art of Problem Solving, 2024a), AMC23/24 (Art
 348 of Problem Solving, 2024b), Minerva Math (Lewkowycz et al., 2022), Gaokao (Zhang et al., 2023),
 349 Olympiad Bench (He et al., 2024). Similar to (Wang et al., 2025b; Zheng et al., 2025b), we eval-
 350 uate models on those benchmarks every 50 steps and report the performance of the checkpoint that
 351 obtains the best average performance on eight benchmarks. For GRPO, we adopt the commonly
 352 used clipping parameter $\epsilon = 0.2$, while for the other baselines, we follow the recommended values
 353 reported in their respective papers. We include more detailed experimental settings in Appendix B.

354 **6.2 PERFORMANCE COMPARISON ON TRAINING WITH STALENESS**
 355

356 **Prosperity without collapse: Stable off-policy training without performance degradation using**
 357 **M2PO.** To verify the effectiveness of M2PO, Table 1 presents a comprehensive comparison of
 358 math reasoning performance across eight benchmarks using models from four different families and
 359 scales, ranging from 1.7B to 32B parameters. We evaluate multiple reinforcement learning methods
 360 under both on-policy and off-policy settings, including GRPO, GSPO, and our proposed M2PO.

361 The results show that while both GRPO and GSPO often suf-
 362 fer significant performance drops under large staleness, M2PO
 363 consistently achieves comparable accuracy to the on-policy
 364 baseline in all training settings. Surprisingly, we notice that,
 365 in some model settings, M2PO with $s = 256$ even achieves a
 366 better performance than M2PO with $s = 0$. For instance, on the
 367 Qwen3-Base-1.7B model, we observe that M2PO with $s = 256$
 368 (36.6%) outperforms GRPO with $s = 0$ (33.0%). A potential
 369 explanation is that small effective staleness (e.g., $s = 0$ corre-
 370 sponding to delays between 0 and 3) can still adversely affect
 371 training stability. Our further analysis in Figure 7 supports this
 372 view, showing that M2PO with $s = 256$ exhibits an even lower
 373 clipping ratio than GRPO with $s = 0$. Overall, these results
 374 show that M2PO remains robust and effective, sustaining stable, high performance even under ex-
 375 treme off-policy conditions.

376 In addition to the final accuracy comparison in Table 1, we also analyze the training dynamics of
 377 accuracy and reward of Qwen-2.5-32B models. As shown in Figure 1, M2PO with $s = 256$ initially
 falls behind the on-policy baseline but quickly catches up, eventually matching its performance,
 while converging much faster and achieving higher accuracy than GRPO under the same staleness.



378 Figure 5: Training reward on
 379 Qwen-2.5-32B.

378
 379 Table 1: Performance (%) comparison across eight math reasoning benchmarks using models from
 380 1.7B to 32B parameters. We report results for GRPO, GSPO, and M2PO under both on-policy
 381 ($s = 0$) and off-policy ($s = 256$) settings. Underlined numbers denote the best average accuracy,
 382 while **bold** numbers highlight the best average accuracy under stale rollouts ($s = 256$). M2PO
 383 consistently improves stability under staleness and achieves higher average accuracy than GRPO.

Method	S	AIME24/25	AMC23/24	Math500	Gaokao	Miner.	Olymp.	Avg.
<i>Llama-3.2-3B-Instruct</i>								
GRPO	0	11.0 / 2.3	31.3 / 17.2	53.6	42.99	23.1	20.3	25.2
GRPO	256	9.6 / 0.4	25.0 / 13.9	52.4	42.08	17.6	18.8	22.5
GSPO	256	9.0 / 0.2	30.0 / 14.4	50.6	40.65	18.8	17.3	22.6
M2PO (Ours)	256	10.4 / 4.4	33.8 / 17.8	52.0	44.48	21.2	18.1	25.3
<i>Qwen2.5-Math-7B</i>								
GRPO	0	39.6 / 17.5	63.8 / 46.7	82.3	64.1	36.7	43.6	49.3
GRPO	256	29.4 / 12.9	64.4 / 39.4	80.5	63.2	33.1	43.1	45.7
GSPO	256	27.3 / 13.1	63.8 / 36.7	79.0	62.2	33.5	41.9	44.7
M2PO (Ours)	256	33.3 / 17.5	63.8 / 40.6	84.0	66.4	38.1	47.1	48.8
<i>Qwen3-Base-1.7B</i>								
GRPO	0	7.5 / 7.5	40.6 / 26.1	67.2	55.9	28.9	30.5	33.0
GRPO	256	8.5 / 4.8	34.4 / 25.0	64.3	52.7	26.0	27.6	30.4
GSPO	256	6.9 / 4.0	39.4 / 18.9	65.0	53.1	26.5	27.5	30.1
M2PO (Ours)	256	14.0 / 6.5	48.1 / 27.8	71.8	59.5	29.4	35.6	36.6
<i>Qwen3-Base-4B</i>								
GRPO	0	22.9 / 20.2	63.8 / 53.9	84.6	69.8	40.2	50.5	50.7
GRPO	256	14.0 / 9.6	51.9 / 32.8	76.8	61.7	34.4	39.8	40.1
GSPO	256	17.9 / 15.4	55.6 / 38.3	76.8	62.3	35.1	44.3	43.2
M2PO (Ours)	256	26.7 / 21.0	64.4 / 49.4	85.8	70.3	40.5	52.3	51.3
<i>Qwen3-Base-8B</i>								
GRPO	0	26.7 / 19.4	76.9 / 52.8	87.7	71.6	41.2	52.8	53.6
GRPO	256	21.0 / 13.1	63.8 / 40.0	81.8	67.8	38.5	47.4	46.7
M2PO (Ours)	256	30.2 / 23.1	71.3 / 56.7	87.2	75.1	42.6	54.8	55.1
<i>Qwen2.5-32B</i>								
GRPO	0	24.4 / 18.3	71.9 / 46.7	85.4	71.9	41.4	52.9	51.6
GRPO	256	20.4 / 9.6	68.1 / 41.1	83.0	67.3	40.9	45.9	47.0
M2PO (Ours)	256	24.8 / 19.4	76.3 / 50.0	85.7	71.7	41.5	51.7	52.6

414 This highlights that M2PO not only maintains comparable final
 415 accuracy but also accelerates convergence when training with
 416 stale data. Figure 5 shows a similar trend in the reward curves.
 417 M2PO with $s = 256$ also starts off behind the on-policy base-
 418 line due to the initial plateau caused by using data generated
 419 from the base model, but it quickly catches up and aligns closely
 420 with the $s = 0$ trajectory. In contrast, GRPO with $s = 256$ con-
 421 sistently underperforms across the entire training trajectory.

422 **Performance of other baselines under staleness.** A number
 423 of prior works have proposed alternative trust region strate-
 424 gies beyond ϵ -clipping, including GSPO (Zheng et al., 2025a),
 425 AREAL (Fu et al., 2025), TOPR (Roux et al., 2025), GPPO (Su
 426 et al., 2025), and CISPO (Chen et al., 2025). *Despite not being*
 427 *designed to handle the extreme staleness studied in this paper*,
 428 we also evaluate these methods under our setting with $s = 256$ for completeness and comparison.
 429 Among these methods, GSPO is the only one that preserves training stability, though it still exhibits
 430 a noticeable performance drop under high staleness (see Table 1). For the other methods, as shown
 431 in Figure 6, most encounter substantial difficulties in maintaining training stability and tend to break
 down early in training. These observations suggest that while existing approaches can be effective

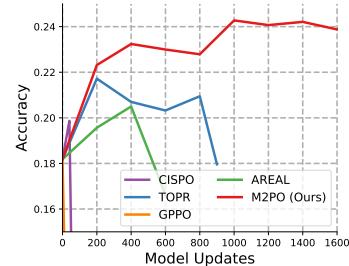


Figure 6: Methods comparison under staleness ($s = 256$) on Llama3.2-Instruct-3B.

under moderately stale settings, they face significant challenges when extended to larger staleness, highlighting the need for a more robust and effective solution.

6.3 ANALYSIS AND ABLATION STUDY

Stable training with reduced clipping.

Figure 7(a)(b) illustrates the clipping dynamics of GRPO and our proposed M2PO under different staleness settings. In Figure 7a, we report results on Qwen-3-Base-1.7B. Under large staleness ($s = 256$), GRPO exhibits frequent clipping events, with the ratio increasing sharply and remaining high during most of the training. In contrast, M2PO under the same staleness maintains an exceptionally low clipping ratio, comparable to or even lower than the on-policy GRPO baseline ($s = 0$). Notably, M2PO with $s = 256$ exhibits less clipping than GRPO with $s = 0$, which explains why M2PO with $s = 256$ achieves higher accuracy than GRPO with $s = 0$ in Table 1. Figure 1 shows the same comparison on Qwen-2.5-32B. A similar trend holds: GRPO with $s = 256$ suffers from substantial clipping, whereas M2PO effectively suppresses unnecessary clipping, remaining close to the on-policy baseline.

Figure 7b summarizes the average clipping ratio across the entire training process. On Qwen-3-Base-1.7B, GRPO with $s = 256$ reaches an average clipping ratio of 0.66%, compared to 0.07% for GRPO with $s = 0$ and only 0.02% for M2PO with $s = 256$. On Qwen-2.5-32B, GRPO with $s = 256$ averages 1.22%, while GRPO with $s = 0$ records 0.05% and M2PO with $s = 256$ maintains a similarly low ratio of 0.06%. These results show that M2PO reduces clipping by over an order of magnitude compared to GRPO, thereby enabling stable and efficient training by clipping only when necessary.

Robustness to the choice of τ_{M_2} . Figure 8 shows that performance is not highly sensitive to the choice of τ_{M_2} . Accuracy remains stable across a broad range, only dropping when τ_{M_2} is set extremely small (overly restrictive constraint) or very large (training collapse). This explains why a single setting of $\tau_{M_2} = 0.04$ works robustly across all training settings in our paper.

Training dynamics on KL and M_2 . Figure 9 shows the impact of M2PO masking on training stability. Figure 9a shows that the average M_2 without masking exhibits frequent spikes throughout training, indicating instability in the second-moment estimates (blue curve). Applying M2PO masking effectively suppresses these fluctuations and maintains consistently low M_2 values, leading to more stable updates (red curve). Figure 9b compares the KL divergence across different methods. Although M2PO with $s = 256$ involves substantially less clipping than GRPO (shown in Figure 7), it maintains a more stable divergence than GRPO with $s = 256$. These results indicate that M2PO performs clipping in a more precise and adaptive manner, ensuring training stability with

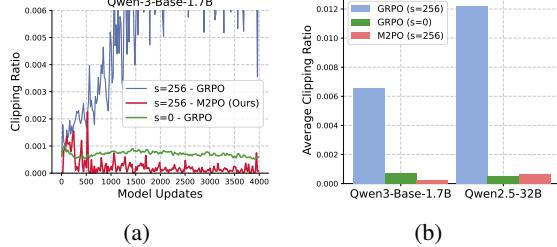


Figure 7: (a) Clipping ratio dynamics during RL on the Qwen-3-Base-1.7B model. (b) Comparison of the average clipping ratio across models and methods.

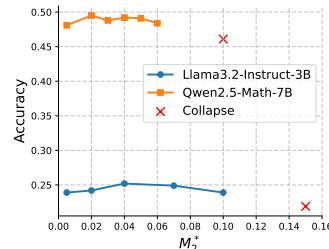


Figure 8: Ablation study of the τ_{M_2} threshold on Llama-3.2-3B-Instruct and Qwen2.5-Math-7B.

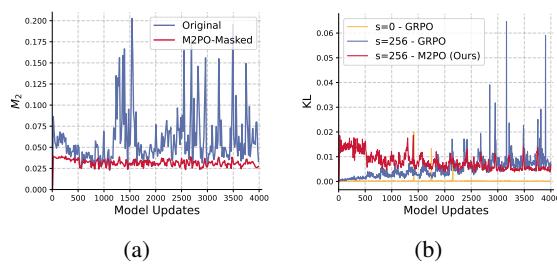


Figure 9: (a) Average M_2 in each model updates with and without M2PO masking on Qwen-2.5-32B, showing that masking effectively suppresses spikes and stabilizes the M_2 throughout training. (b) Average KL divergence in each model updates on Qwen-2.5-32B under different methods.

486 substantially less reliance on clipping. These results demonstrate that M2PO enables more pre-
 487 cise and adaptive clipping, achieving training stability while relying on significantly fewer clipping
 488 operations, and thereby attaining better performance without risking training collapse.
 489

490 7 CONCLUSION

491
 492 In this work, we investigated why off-policy RL for LLMs often fails under stale data and uncov-
 493 ered the *prosperity-before-collapse* phenomenon: training without a trust region initially outper-
 494 forms standard methods, showing that stale data can be as informative as on-policy trajectories,
 495 but eventually collapses due to instability. Motivated by this observation, we proposed **M2PO**,
 496 which constrains the second moment of importance weights to provide a variance-sensitive and sta-
 497 ble trust region. This design suppresses extreme outliers while preserving informative high-entropy
 498 tokens, enabling stable training that matches on-policy performance even under extreme staleness.
 499 Extensive experiments further demonstrate that M2PO significantly reduces clipping and is highly
 500 insensitive to its threshold, highlighting its practicality and scalability for efficient RL with LLMs.
 501

502 8 ETHICS STATEMENT

503
 504 This work focuses on developing reinforcement learning algorithms for large language models. Our
 505 research does not involve human subjects, personally identifiable information, or sensitive data.
 506 All datasets used are publicly available and widely adopted in the community. We acknowledge that
 507 more capable LLMs may have potential societal impacts, including misuse for generating misleading
 508 or harmful content. To mitigate these risks, our study is confined to controlled academic settings,
 509 and our primary goal is to improve the stability and efficiency of training methods.
 510

511 9 REPRODUCIBILITY STATEMENT

512
 513 All implementation details, hyperparameters, and experimental setups are thoroughly documented
 514 in the paper and appendix, providing sufficient information for independent reproduction of our
 515 results.
 516

517 REFERENCES

518
 519 Charles Arnal, GaĂłtan Narozniak, Vivien Cabannes, Yunhao Tang, Julia Kempe, and Remi
 520 Munos. Asymmetric reinforce for off-policy reinforcement learning: Balancing positive and neg-
 521 ative rewards. *arXiv preprint arXiv:2506.20520*, 2025.

522 Art of Problem Solving. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions,
 523 2024a. Accessed: 2025-04-20.

524 Art of Problem Solving. Amc problems and solutions. https://artofproblemsolving.com/wiki/index.php?title=AMC_Problems_and_Solutions, 2024b. Accessed: 2025-04-20.

525
 526 Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu,
 527 Chao Wang, Cheng Zhu, et al. Minimax-m1: Scaling test-time compute efficiently with lightning
 528 attention. *arXiv preprint arXiv:2506.13585*, 2025.

529 Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen
 530 Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for
 531 reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.

532 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
 533 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
 534 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
 535 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,

540 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
 541 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
 542 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai
 543 Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,
 544 Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang,
 545 Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang,
 546 Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang,
 547 R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhua Chen, Shengfeng
 548 Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing
 549 Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen
 550 Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong
 551 Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu,
 552 Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xi-
 553 aoshua Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia
 554 Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng
 555 Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong
 556 Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong,
 557 Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou,
 558 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying
 559 Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda
 560 Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu,
 561 Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu
 562 Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforce-
 563 ment learning. *arXiv preprint arXiv:2501.12948*, 2025.

564 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 565 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 566 *arXiv e-prints*, pp. arXiv–2407, 2024.

567 Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun
 568 Mei, Jiashu Wang, et al. Areal: A large-scale asynchronous reinforcement learning system for
 569 language reasoning. *arXiv preprint arXiv:2505.24298*, 2025.

570 Jiaxuan Gao, Shusheng Xu, Wenjie Ye, Weilin Liu, Chuyi He, Wei Fu, Zhiyu Mei, Guangju Wang,
 571 and Yi Wu. On designing effective rl reward at training time for llm reasoning. *arXiv preprint
 572 arXiv:2410.15115*, 2024.

573 Zitian Gao, Lynx Chen, Haoming Luo, Joey Zhou, and Bryan Dai. One-shot entropy minimization.
 574 *arXiv preprint arXiv:2505.20282*, 2025.

575 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,
 576 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for
 577 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint
 578 arXiv:2402.14008*, 2024.

579 Jingkai He, Tianjian Li, Erhu Feng, Dong Du, Qian Liu, Tao Liu, Yubin Xia, and Haibo
 580 Chen. History rhymes: Accelerating llm reinforcement learning with rhymerl. *arXiv preprint
 581 arXiv:2508.18588*, 2025.

582 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
 583 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
 584 2021.

585 Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv
 586 preprint arXiv:2501.03262*, 2025.

587 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
 588 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
 589 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

590 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
 591 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint
 592 arXiv:2310.06770*, 2023.

594 Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy,
 595 Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through
 596 refined credit assignment. *arXiv preprint arXiv:2410.01679*, 2024.

597

598 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
 599 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
 600 serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023.

601

602 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
 603 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
 604 reasoning problems with language models. *Advances in Neural Information Processing Systems*,
 605 35:3843–3857, 2022.

606

607 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittweis, Rémi Leblond, Tom
 608 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien
 609 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven
 610 Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Push-
 611 meet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code
 612 generation with alphacode. *Science*, 378(6624):1092–1097, 2022. doi: 10.1126/science.abq1158.
 613 URL <https://www.science.org/doi/abs/10.1126/science.abq1158>.

614

615 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
 616 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
 617 *arXiv:2305.20050*, 2023.

618

619 Liyuan Liu, Feng Yao, Dinghuai Zhang, Chengyu Dong, Jingbo Shang, and Jianfeng Gao. Flashrl:
 620 8bit rollouts, full power rl, August 2025a. URL <https://fengyao.notion.site/flash-rl>.

621

622 Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee,
 623 and Min Lin. Understanding rl-zero-like training: A critical perspective. *arXiv preprint*
 624 *arXiv:2503.20783*, 2025b.

625

626 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-
 627 ence on Learning Representations*, 2019.

628

629 Michael Luo, Sijun Tan, Roy Huang, Xiaoxiang Shi, Rachel Xin, Colin Cai, Ameen Patel, Alpay
 630 Ariyak, Qingyang Wu, Ce Zhang, et al. Deepcoder: A fully open-source 14b coder at o3-mini
 631 level, 2025a.

632

633 Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin
 634 Cai, Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler:
 635 Surpassing o1-preview with a 1.5b model by scaling rl. <https://pretty-radio-b75.notion.site/DeepScalE-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>, 2025b. Notion
 636 Blog.

637

638 Michael Noukhovitch, Shengyi Huang, Sophie Xhonneux, Arian Hosseini, Rishabh Agarwal, and
 639 Aaron Courville. Asynchronous rlhf: Faster and more efficient off-policy rl for language models.
 640 *arXiv preprint arXiv:2410.18252*, 2024.

641

642 OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden
 643 Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko,
 644 Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally
 645 Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich,
 646 Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghor-
 647 bani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botaao Hao,
 648 Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary
 649 Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang,
 650 Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel
 651 Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson,

648 Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Eliz-
 649 abeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang,
 650 Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred
 651 von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace
 652 Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart An-
 653 drin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan,
 654 Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever,
 655 Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng,
 656 Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish,
 657 Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan
 658 Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl
 659 Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu,
 660 Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam
 661 Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kon-
 662 draciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen,
 663 Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet
 664 Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael
 665 Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles
 666 Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil
 667 Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg
 668 Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov,
 669 Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar
 670 Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan
 671 Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agar-
 672 wal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu,
 673 Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph
 674 Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Tay-
 675 lor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson,
 676 Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna
 677 Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi
 678 Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen,
 Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li.
 679 Openai o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.

680 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
 681 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
 682 low instructions with human feedback. *Advances in neural information processing systems*, 35:
 683 27730–27744, 2022.

684 Nicolas Le Roux, Marc G Bellemare, Jonathan Lebensold, Arnaud Bergeron, Joshua Greaves,
 685 Alex Fréchette, Carolyne Pelletier, Eric Thibodeau-Laufer, Sándor Toth, and Sam Work. Ta-
 686 pered off-policy reinforce: Stable and efficient reinforcement learning for llms. *arXiv preprint*
 687 *arXiv:2503.14286*, 2025.

688 John Schulman. The KL divergence between two close distributions. [http://joschu.net/](http://joschu.net/blog/kl-approx.html)
 689 blog/kl-approx.html, 2015. Accessed: 2025-11-18.

690 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
 691 policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR,
 692 2015.

693 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
 694 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

695 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
 696 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical
 697 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

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

702 Zhenpeng Su, Leiyu Pan, Xue Bai, Dening Liu, Guanting Dong, Jiaming Huang, Wenping Hu, and
 703 Guorui Zhou. Klear-reasoner: Advancing reasoning capability via gradient-preserving clipping
 704 policy optimization. *arXiv preprint arXiv:2508.07629*, 2025.

705 Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun
 706 Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with
 707 llms. *arXiv preprint arXiv:2501.12599*, 2025.

708 Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,
 709 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive
 710 effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025a.

711 Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
 712 Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software
 713 developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.

714 Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Lucas Liu, Baolin Peng, Hao Cheng, Xuehai
 715 He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language
 716 models with one training example. *arXiv preprint arXiv:2504.20571*, 2025b.

717 Yongji Wu, Xueshen Liu, Haizhong Zheng, Juncheng Gu, Beidi Chen, Z Morley Mao, Arvind
 718 Krishnamurthy, and Ion Stoica. Rlboost: Harvesting preemptible resources for cost-efficient re-
 719 inforcement learning on llms. *arXiv preprint arXiv:2510.19225*, 2025.

720 Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong
 721 Zhang, Caiming Xiong, et al. A minimalist approach to llm reasoning: from rejection sampling
 722 to reinforce. *arXiv preprint arXiv:2504.11343*, 2025.

723 Ran Yan, Youhe Jiang, Tianyuan Wu, Jiaxuan Gao, Zhiyu Mei, Wei Fu, Haohui Mai, Wei Wang,
 724 Yi Wu, and Binhang Yuan. Areal-hex: Accommodating asynchronous rl training over heteroge-
 725 neous gpus. *arXiv preprint arXiv:2511.00796*, 2025.

726 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-
 727 hong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical
 728 expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.

729 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
 730 Chang Gao, Chengan Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint
 731 arXiv:2505.09388*, 2025.

732 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong
 733 Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at
 734 scale. *arXiv preprint arXiv:2503.14476*, 2025.

735 Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang,
 736 TianTian Fan, Zhengyin Du, Xiangpeng Wei, et al. Vapo: Efficient and reliable reinforcement
 737 learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*, 2025a.

738 Yufeng Yuan, Yu Yue, Ruofei Zhu, Tiantian Fan, and Lin Yan. What's behind ppo's collapse in
 739 long-cot? value optimization holds the secret. *arXiv preprint arXiv:2503.01491*, 2025b.

740 Xiaojiang Zhang, Jinghui Wang, Zifei Cheng, Wenhao Zhuang, Zheng Lin, Minglei Zhang, Shao-
 741 jie Wang, Yinghan Cui, Chao Wang, Junyi Peng, Shimiao Jiang, Shiqi Kuang, Shouyu Yin,
 742 Chaohang Wen, Haotian Zhang, Bin Chen, and Bing Yu. Srpo: A cross-domain implemen-
 743 tation of large-scale reinforcement learning on llm, 2025. URL <https://arxiv.org/abs/2504.14286>.

744 Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying, Liang He, and Xipeng Qiu. Evaluating the
 745 performance of large language models on gaokao benchmark. *arXiv preprint arXiv:2305.12474*,
 746 2023.

747 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright,
 748 Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully
 749 sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023.

756 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang,
757 Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint*
758 *arXiv:2507.18071*, 2025a.

759
760 Haizhong Zheng, Yang Zhou, Brian R Bartoldson, Bhavya Kailkhura, Fan Lai, Jiawei Zhao, and
761 Beidi Chen. Act only when it pays: Efficient reinforcement learning for llm reasoning via selective
762 rollouts. *arXiv preprint arXiv:2506.02177*, 2025b.

763 Yinmin Zhong, Zili Zhang, Xiaoniu Song, Hanpeng Hu, Chao Jin, Bingyang Wu, Nuo Chen, Yukun
764 Chen, Yu Zhou, Changyi Wan, et al. Streamrl: Scalable, heterogeneous, and elastic rl for llms
765 with disaggregated stream generation. *arXiv preprint arXiv:2504.15930*, 2025.

766
767 Zilin Zhu, Chengxing Xie, Xin Lv, and slime Contributors. slime: An llm post-training framework
768 for rl scaling. <https://github.com/THUDM/slime>, 2025. GitHub repository. Corresponding
769 author: Xin Lv.

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810 **A OVERVIEW**
811812 In this appendix, we provide additional details to complement the main text. Appendix B de-
813 scribes the experimental setup in full, including datasets, models, and training hyperparameters.
814 Appendix D presents theoretical proofs supporting our method and analysis.
815816 **B DETAILED EXPERIMENTAL SETTING**
817818 **Models & Datasets.** To verify the effectiveness of our method, we extensively evaluate M2PO
819 on six models from four model series: Qwen2.5-Math-7B (Yang et al., 2024), Llama-3.2-3B-
820 Instruct (Dubey et al., 2024), Qwen3-Base-1.7B/4B/8B (Yang et al., 2025), and Qwen2.5-32B (Yang
821 et al., 2024). For Qwen2.5-Math-7B, we use a context length of 4k, which is the maximum for this
822 series. while for all other models the context length is set to 16k. For training, we adopt the Deep-
823 ScaleR (Luo et al., 2025b) math dataset.
824825 **Training.** Our method is implemented based on verl (Sheng et al., 2024) pipeline and uses
826 vLLM (Kwon et al., 2023) for rollout. We use a mix of H100 and H200 servers for training, de-
827 pending on resource availability. We set the rollout temperature to 1 for vLLM (Kwon et al., 2023).
828 The training batch size is set to 256 prompts, and the mini-batch size to 64 prompts. We sample 8
829 responses per prompt. We train all models for 1000 steps, and we optimize the actor model using
830 the AdamW (Loshchilov & Hutter, 2019) optimizer with a constant learning rate of 1e-6. We use
831 $\beta_1 = 0.9$, $\beta_2 = 0.999$, and apply a weight decay of 0.01. We use the following question template to
832 prompt the LLM. For reward assignment, we give a score of 0.1 for successfully extracting an an-
833 swer and a score of 1.0 if the extracted answer is correct. Similar to (Yu et al., 2025), we remove the
834 KL-divergence term. The optimization is performed on the parameters of the actor module wrapped
835 with Fully Sharded Data Parallel (FSDP) (Zhao et al., 2023) for efficient distributed training. We set
the M2PO threshold to 0.04 for all training runs.
836837 **Evaluation.** For benchmark datasets, we use eight widely used complex mathematical reason-
838 ing benchmarks to evaluate the performance of trained models: Math500 (Hendrycks et al., 2021;
839 Lightman et al., 2023), AIME24/25 (Art of Problem Solving, 2024a), AMC23/24 (Art of Problem
840 Solving, 2024b), Minerva Math (Lewkowycz et al., 2022), Gaokao (Zhang et al., 2023), Olympiad
841 Bench (He et al., 2024). Same as the training setting, For Qwen2.5-Math-7B models, we use 4k
842 as the context length. For other models, we set the context length to 16k. Similar to (Wang et al.,
843 2025b; Zheng et al., 2025b), we evaluate models on those benchmarks every 50 steps and report the
844 performance of the checkpoint that obtains the best average performance on eight benchmarks. We
845 evaluate all models with temperature = 1. For AIME24/25, we report the *pass@1(avg@16)*, for
other benchmarks, we report the *pass@1(avg@4)*.
846847 **Baselines.** For all baseline methods, we adopt exactly the same training setup and evaluation proto-
848 col to ensure a fair comparison. For trust-region settings, we follow the recommended hyperparam-
849 eters reported in the original papers whenever available. Specifically, for GSPO, we use $3e-4$ and
850 $4e-4$ as the lower and upper bounds. For CISPO, we use a clipping threshold of 0.2, as the original
851 paper does not specify a particular value; we also experimented with a looser bound, but found that
852 it makes training significantly more prone to collapse. For AREAL, we use a clipping threshold of
0.2. For GPPO, we set the lower and upper clipping bounds to 0.2 and 0.28, respectively. TOPR
853 does not contain any trust-region hyperparameters.
854855 **Question Template**856 Please solve the following math problem: {{Question Description}}. The assistant first thinks
857 about the reasoning process step by step and then provides the user with the answer. Return
858 the final answer in \boxed{} tags, for example \boxed{1}. Let's solve this step by step.
859860 **C THE USE OF LARGE LANGUAGE MODELS (LLMs)**
861862 We used large language models (LLMs) as general-purpose assistants in two limited ways: (1) for
863 writing polish, including improving grammar, readability, and presentation of the manuscript, and
864

(2) as code assistants (e.g., Cursor, GitHub Copilot) to accelerate routine coding tasks such as debugging syntax errors and refactoring simple functions. LLMs were not used for research ideation, algorithm design, experimental analysis, or drawing conclusions. All conceptual and scientific contributions are entirely the work of the authors.

D THEORETICAL PROOF

Proof of Theorem 1.

Proof. Let $z = \log r$, so that $r = e^z$. Since $1/R < r \leq R$, we have $|z| \leq \log R$.

For $z \geq 0$, observe that

$$\frac{e^z - 1}{z} = \int_0^1 e^{tz} dt \leq \int_0^1 e^z dt = e^z,$$

which implies $(e^z - 1)^2 \leq (ze^z)^2 = z^2 e^{2z}$.

For $z \leq 0$, set $u = -z \geq 0$. Then

$$(e^z - 1)^2 = (1 - e^{-u})^2 \leq u^2 = z^2 \leq z^2 e^{2|z|}.$$

Combining both cases, for all $z \in \mathbb{R}$ we obtain

$$(e^z - 1)^2 \leq z^2 e^{2|z|}.$$

Substituting $Z = \log r$, this yields

$$(r - 1)^2 \leq (\log r)^2 e^{2|\log r|} \leq R^2 (\log r)^2.$$

Taking expectation under π_{behav} gives

$$\chi^2(\pi_{\text{new}} \parallel \pi_{\text{behav}}) = \mathbb{E}_{\pi_{\text{behav}}}[(r - 1)^2] \leq R^2 \mathbb{E}_{\pi_{\text{behav}}}[(\log r)^2] = R^2 M_2.$$

□

E THEORETICAL INSIGHTS ON PROSPERITY BEFORE COLLAPSE AND M2PO

In this section, we analyze why the unclipped objective is unbiased, while the ε -clipped PPO objective introduces bias by construction. This explains the observed *prosperity-before-collapse* phenomenon.

Unclipped gradient is unbiased. Let the importance ratio be

$$r_t(\theta) = \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\text{behav}}(a_t \mid s_t)}.$$

The true policy gradient under off-policy sampling from π_{behav} is

$$\nabla_\theta J(\theta) = \mathbb{E}_{t \sim \pi_{\text{behav}}} [r_t(\theta) A_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)].$$

The surrogate objective without clipping is

$$L_{\text{unclip}}(\theta) = \mathbb{E}_t [r_t(\theta) A_t],$$

whose gradient matches the true policy gradient,

$$\nabla_\theta L_{\text{unclip}}(\theta) = \mathbb{E}_t [r_t(\theta) A_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)],$$

and is therefore strictly unbiased:

$$\mathbb{E}[\nabla_\theta L_{\text{unclip}}(\theta)] = \nabla_\theta J(\theta).$$

PPO ε -clipping introduces bias. The clipped surrogate is

$$L_{\text{clip}}(\theta) = \mathbb{E}_t [\min(r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon) A_t)].$$

918 Consider the set of samples where $(A_t > 0, r_t(\theta) > 1 + \varepsilon)$. For these samples PPO replaces $r_t(\theta)$
 919 with a constant $(1 + \varepsilon)$:

$$920 \quad \min(r_t A_t, (1 + \varepsilon) A_t) = (1 + \varepsilon) A_t,$$

921 whose gradient is zero:

$$922 \quad \nabla_{\theta} [(1 + \varepsilon) A_t] = 0.$$

923 However, the true gradient contribution from these same samples is

$$925 \quad \nabla_{\theta} (r_t A_t) = r_t(\theta) A_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \neq 0.$$

926 Thus the gradient contribution over this region is removed. Let $\Omega_+ \cup \Omega_-$ denote the set of clipped
 927 samples. Then,

$$928 \quad \nabla_{\theta} L_{\text{clip}}(\theta) = \mathbb{E}_{t \notin (\Omega_+ \cup \Omega_-)} [r_t A_t \nabla_{\theta} \log \pi_{\theta}],$$

929 whereas the true gradient is

$$930 \quad \nabla_{\theta} J(\theta) = \mathbb{E}_t [r_t A_t \nabla_{\theta} \log \pi_{\theta}].$$

932 We denote by Ω_+ and Ω_- the sets of timesteps where ε -clipping becomes active:

$$933 \quad \Omega_+ = \{t \mid A_t > 0, r_t(\theta) > 1 + \varepsilon\},$$

$$935 \quad \Omega_- = \{t \mid A_t < 0, r_t(\theta) < 1 - \varepsilon\}.$$

936 Then the difference introduced by clipping is

$$938 \quad \Delta g(\theta) = -\mathbb{E}_{t \in (\Omega_+ \cup \Omega_-)} [r_t(\theta) A_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] \neq 0,$$

939 which shows that PPO's clipped gradient is *biased by construction* whenever clipping occurs.

940 Combining these properties yields the observed **prosperity-before-collapse phenomenon**: When
 941 the clipping is removed, the surrogate gradient becomes the unbiased policy gradient, enabling faster
 942 early-stage improvement. However, without a trust region, the importance ratios under off-policy
 943 rollouts quickly become heavy-tailed, causing the gradient variance to grow rapidly, eventually lead-
 944 ing to training collapse.

945 The reason M2PO achieves a better trade-off between performance and stability is that its masking
 946 strategy is more adaptive. By using the second moment to quantify batch-level distribution shift,
 947 M2PO automatically adjusts the effective trust region: when variance is small, it keeps the clipping
 948 region wide; when variance grows, it tightens the region accordingly. As a result, the measure of
 949 clipped samples, $|\Omega_+ \cup \Omega_-|$, is significantly reduced compared to fixed ε -clipping, leading to a
 950 smaller bias while still preventing variance explosion.

952 F DISCUSSION AND ANALYSIS

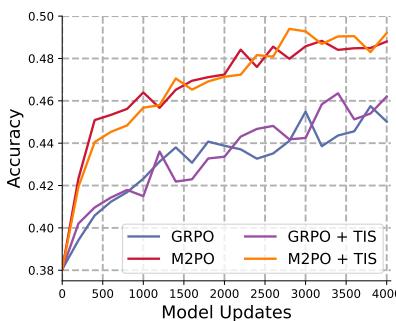
954 F.1 WHY DOES TOLERANCE TO HIGH STALENESS MATTER FOR RL ON LLMs?

956 High staleness tolerance is both practical and increasingly unavoidable for future RL for LLMs and
 957 for scalable RL system designs. Here are two key reasons:

959 1. High staleness tolerance enables training on more complex tasks. Complex tasks can substantially
 960 increase rollout latency. This is especially true for coding and agent environments, where **tool**
 961 **calls** and **reward computation** can be expensive, as both may require executing external programs,
 962 interacting with third-party APIs, or running large numbers of test cases to verify correctness.

963 In realistic settings, such as SWE-bench Jimenez et al. (2023) with OpenHands Wang et al. (2024),
 964 rollout latency can become extremely large. A single run, even with a small batch size, the end-
 965 to-end inference time (including tool call and code execution) can exceed 100 minutes, with more
 966 than 80 iterations with the environment. For reward calculation, verifying correctness against tens
 967 or hundreds of test cases also requires repeatedly compiling and running programs, adding another
 968 5–20 minutes of reward computation depending on code and test-case complexity.

969 As a result, the combined rollout + reward latency can easily reach 120 minutes or more. In contrast,
 970 policy updates in LLM RL, especially with sufficient parallelism, can be as fast as 60 seconds per
 971 update. Even without any additional system-induced delays, this already corresponds to up to 120
 972 model updates of staleness, and the staleness can be even larger for more complex tasks.



972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
Figure 10: Combining TIS with GRPO and M2PO on Qwen-2.5-Math-7B with $s = 256$. Combined with TIS, M2PO shows a slight performance improvement, as TIS better mitigates the distribution gap between FSDP and VLLM. However, TIS alone cannot address the off-policy caused by staleness.

To enable future RL to support increasingly complex LLM agents, especially those involving frequent tool calls and substantial code execution, training under large staleness is practically necessary.

2. High staleness tolerance enables flexible and scalable system design. Modern LLM RL training systems increasingly adopt decoupled and asynchronous designs to improve throughput and reduce cost. In such architectures, the ability to tolerate high staleness directly enables more flexible and scalable RL pipelines for LLMs Fu et al. (2025); Zhu et al. (2025); Wu et al. (2025); Yan et al. (2025), like *heterogeneous computing and disaggregated computing systems*.

1) **Heterogeneous Computing.** As demonstrated in a recent work Yan et al. (2025), offloading rollouts to slower but cost-efficient machines can substantially reduce the overall training expense. However, this will significantly increase the rollout latency, leading to larger staleness.

2) **Disaggregated Computing.** In large-scale RL training systems, rollout generation and policy updates can run on physically separated clusters or even across different data centers. Network bandwidth and latency can become a huge bottleneck: synchronizing model weights and transferring rollout trajectories can further amplify staleness as the system scales.

To sum up, these factors make high staleness a realistic and sometimes unavoidable operating regime in modern RL for LLMs. Tolerance to high staleness is therefore crucial not only for cost-efficient and flexible system design, but also for enabling RL on complex long-latency tasks without sacrificing throughput.

Motivated by this, we use a clean and controlled stale-k abstraction to isolate the effect of data staleness on RL algorithms. Within this controlled but representative setup, we show the limitation of GRPO under high staleness, while M2PO remains stable and approaches on-policy performance. We believe this work provides algorithmic understandings of RL under high staleness and lays the foundation for future system-level integration.

F.2 HOW IS M2PO RELATED TO TIS (TRUNCATED IMPORTANCE SAMPLING)?

Another off-policy RL technique, TIS (Liu et al., 2025a), addresses a different source of off-policy drift: the mismatch between the inference engine used for rollouts (e.g., vLLM, SGLang, or quantized engines) and the training engine (e.g., FSDP or Megatron). Despite that the weights are exactly the same, the engine difference on training and inference can introduce sampling differences and cause potential mismatch issues. This form of off-policy discrepancy is orthogonal to the source we study: **data staleness**.

More specifically, the optimization objective for TIS in (Liu et al., 2025a) is:

$$\mathbb{E}_{\pi_{\text{vllm}}(\theta_{\text{behav}})} \left[\min \left(\frac{\pi_{\text{fsdp}}(a, \theta_{\text{behav}})}{\pi_{\text{vllm}}(a, \theta_{\text{behav}})}, C \right) \cdot \nabla_{\theta} \min \left(\frac{\pi_{\text{fsdp}}(a, \theta)}{\pi_{\text{fsdp}}(a, \theta_{\text{behav}})} \hat{A}, \text{clip} \left(\frac{\pi_{\text{fsdp}}(a, \theta)}{\pi_{\text{fsdp}}(a, \theta_{\text{behav}})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A} \right) \right]$$

We can see that the TIS objective introduces a truncated (or scaled) importance ratio to prevent excessively large ratios and stabilize training under inference-engine mismatch. However, for the

1026 off-policy drift caused by staleness, TIS still relies on a fixed ϵ -based clipping region to form a
 1027 trust region. In contrast, the primary goal of M2PO is to provide an adaptive, variance-sensitive
 1028 mechanism to determine the trust region based on the second-moment statistics of importance ratios,
 1029 without relying on a hand-chosen ϵ . In this sense, TIS and M2PO are complementary: TIS mitigates
 1030 inference-engine mismatch; M2PO stabilizes training under stale data and large-degree off-policy
 1031 drift. Here is a formula combining TIS and M2PO:
 1032
 1033

$$\mathcal{J}_{\text{M2PO}}(\theta) = \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min \left(\frac{\pi_{\theta_{\text{behav}}, \text{fsdp}}(o_i|q)}{\pi_{\theta_{\text{behav}}, \text{vllm}}(o_i|q)}, C \right) M_{i,t} \frac{\pi_{\theta, \text{fsdp}}(o_i|q)}{\pi_{\theta_{\text{behav}}, \text{fsdp}}(o_i|q)} A_{i,t}, \quad (6)$$

1037 where $M_{i,t}$ is the mask generated by Algorithm 1.
 1038

1039 In Figure 10, we evaluate the combination of TIS with both GRPO and M2PO on Qwen2.5-Math-
 1040 7B under a high-staleness setting ($s = 256$). When integrated with TIS, M2PO achieves a slight
 1041 performance improvement, as TIS helps reduce the distribution gap between the FSDP training
 1042 engine and the VLLM rollout engine. (That said, the gain is modest because both training and
 1043 rollout are performed in FP16, which keeps the engine-induced distribution shift relatively small.)
 1044 However, TIS alone is insufficient to address the distribution shift introduced by large staleness.
 1045

1046 F.3 COMPARE M2PO WITH ASYMMETRIC CLIPPING

1047 We have conducted an additional experiment using the
 1048 asymmetric clipping ratio (0.2, 0.28) used in DAPO for
 1049 our setting and report the results in Figure 11. Despite
 1050 the asymmetric clipping ratio improving the GRPO per-
 1051 formance, it still underperforms M2PO.
 1052

1053 The fundamental distinction between M2PO and asym-
 1054 metric clipping is twofold:
 1055

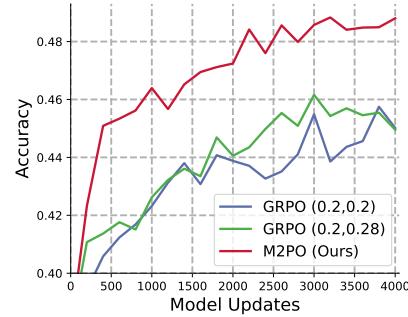
1) **M2PO adaptively determines the trust region based**
 1056 **on the actual distribution shift of each training batch.**
 1057 The second-moment statistic automatically tightens the
 1058 trust region when the shift becomes large, and relaxes it
 1059 when the shift is small, thereby preserving more infor-
 1060 mative gradients. This per-batch adaptivity cannot be re-
 1061 produced by any static ϵ -clipping threshold: a fixed up-
 1062 per/lower bound does not respond to the evolving off-
 1063 policy drift during training.

1064 2) **Tuning clipping ratios is highly brittle and fundamentally impractical for different settings.**
 1065 The optimal asymmetric bounds vary across model sizes, staleness levels, and even across different
 1066 phases of training. The tuning complexity of asymmetric clipping becomes even higher because
 1067 both upper and lower bounds must be selected jointly (quadratic complexity). In contrast, the M2PO
 1068 threshold directly reflects the current batch’s distribution shift, which is the reason that the M2PO
 1069 threshold is not sensitive, and we can use the same threshold in all our experiments.

1070 We believe that these two reasons illustrate why simply adjusting clipping boundaries is not suffi-
 1071 cient for stabilizing stale off-policy training and why the community continues to develop principled
 1072 trust-region algorithms rather than relying on hand-tuned clipping thresholds.
 1073

1074 F.4 M2PO COMPUTATION COST ANALYSIS

1076 Although we use per-token second-moment (M_2) estimates to compute the batch-level M_2 statistic,
 1077 we do not require logits or probabilities over the entire vocabulary. Instead, M_2 is computed
 1078 only from the sampled token at each position (i.e., a single-sample Monte Carlo estimate). Thus,
 1079 M_2 estimation adds no additional forward-pass computation and no additional tensor storage, avoiding
 the $O(\text{VocabularySize})$ cost.



1075 Figure 11: The performance comparison between M2PO and GRPO on coding tasks.

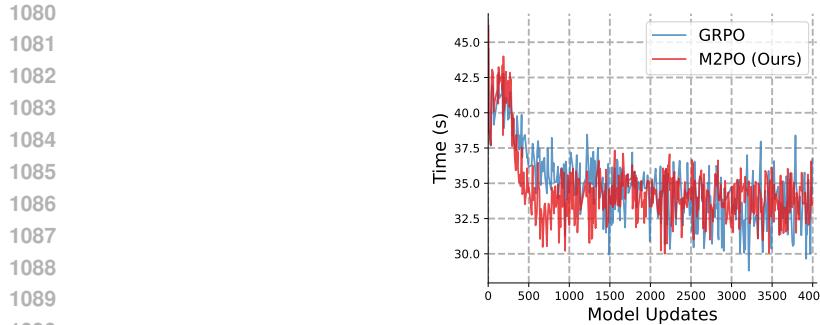


Figure 12: Training time (i.e., update policy) comparison between GRPO and M2PO. M2PO masking does not introduce a noticeable overhead in training stage.

Table 2: Comparison of computation time between GRPO and M2PO. Loss computation contributes a negligible portion of the total training time.

	Loss Compute Time (s)	Other Training Time (s)
GRPO	0.038 (0.11%)	34.86 (99.89%)
M2PO (Ours)	0.065 (0.19%)	34.43 (99.81%)

Besides, for the masking procedure, our implementation sorts tokens by their M_2 values once and iteratively masks those with the largest M_2 until the remaining tokens’ average M_2 falls below the threshold τ_{M_2} . This procedure has a time complexity of $O(N \log N)$, where N is the number of tokens in a batch. In practice, N is typically a few thousand, and this sorting step is negligible compared to the forward/backward passes of large language models. As we illustrated in Figure 12, GRPO and M2PO show a similar training time, which indicates that M2PO masking does not introduce a noticeable overhead in the training stage.

We further provide a breakdown of the loss-computation overhead (including masking computation) during the entire training process. Table 2 below reports the average loss-computation time and the remaining training time (including forward and backward passes) for RL training on Qwen2.5-Math-7B. As shown, M2PO introduces a slightly higher computational cost than GRPO due to the additional sorting required to select tokens with the highest M_2 values. However, this overhead is extremely small relative to the full training step and is far from the bottleneck. For instance, M2PO spends only 0.19% of total time on loss computation and does not slow down the overall training process.

Overall, both M_2 estimation and the M2PO masking procedure incur **minimal computational and memory overhead**, and do not slow down training or increase storage requirements in our experiments.

F.5 HOW IS M2PO RELATED TO KL CONSTRAINT/LOSS?

The KL-divergence loss calculated in the Deepseek R1 (DeepSeek-AI et al., 2025) is a low-variance approximation of the true KL divergence between the old and new policies (Schulman, 2015), which can partially mitigate off-policy drift.

In our early exploration, we also experimented with adding a KL loss. While this regularizer does enhance stability, we found that a trust region is still required to prevent collapse. Once a trust region is introduced, however, we immediately encounter the same challenge discussed in Q3: *how to select an adaptive and effective trust region boundary that balances stability and performance*. Simply tuning a KL penalty or fixed KL target does not adequately address the dynamic distribution shift introduced by stale data.

Moreover, the KL constraint and M2PO are complementary. The KL penalty controls global deviation between the new and old policies, whereas M2PO adaptively masks high-shift tokens based on per-batch second-moment statistics. Combining the two is straightforward and can potentially further improve stability. In this work, we intentionally did not include a KL penalty in our main

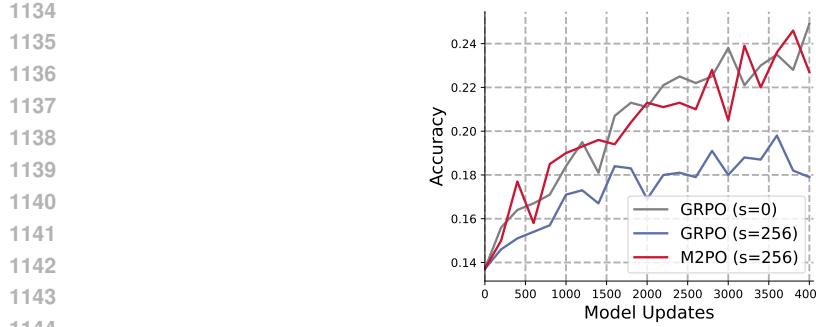


Figure 13: The performance comparison between M2PO and GRPO on coding tasks.

experiments. Our goal is to isolate and understand the role of the trust-region mechanism itself under stale off-policy training. We will explore the combination of KL regularization and M2PO as a promising future direction.

G ADDITIONAL EXPERIMENTS

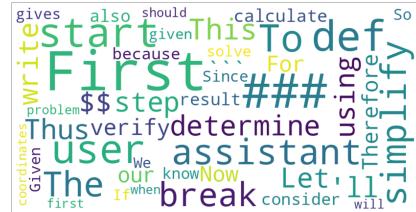
G.1 EVALUATION RESULTS ON CODING TASKS.

We further evaluate M2PO on coding tasks beyond the math domain. Specifically, we train **DeepSeek-R1-Distill-Qwen-1.5B** DeepSeek-AI et al. (2025) on the `code_contests` dataset (Li et al., 2022) and evaluate on `LiveCodeBench` (Jain et al., 2024). We follow a training setup similar to that used for math tasks, except that each training iteration samples a batch of 128 prompts for rollout and 32 prompts for each model updates. We report **avg@4** accuracy as the evaluation metric.

As shown in Figure 13, M2PO with staleness $s = 256$ outperforms GRPO at the same staleness level and achieves performance comparable to GRPO with $s = 0$, consistent with our observations on math reasoning tasks. This further demonstrates the generality of M2PO on other tasks beyond math.

G.2 COMMONLY CLIPPED TOKENS IN GRPO.

Figure 14 shows the specific tokens that are most frequently clipped by ϵ -clipping. The word cloud of commonly clipped tokens is highly aligned with high-entropy tokens shown in (Wang et al., 2025a): these tokens are not random or unimportant, but rather belong to the most semantically and structurally critical elements in reasoning traces. Many of them (e.g., `First`, `simplify`, `determine`, `To def`, `Thus`, `verify`, `break`) are precisely the high-entropy “pivotal tokens” that initiate, connect, or conclude key reasoning steps. Others (e.g., `assistant`, `user`, `code markers` like `###` or `$$`) serve as structural anchors in the dialogue or mathematical formatting. This observation aligns with Figure 4b: as the importance weight ratio $|r - 1|$ grows, clipped tokens tend to exhibit higher entropy.

Figure 14: Word clouds of frequently clipped tokens with ϵ clipping.

G.3 M2PO UNDER DIFFERENT STALENESS.

In Table 1, we compare the performance between M2PO and other baselines under $s = 256$. To further study the robustness of M2PO across different staleness levels, we train Qwen2.5-Math-7B with M2PO under $s = 0, 32, 64, 128, 256$, as shown in Figure 15. We observe that larger staleness slightly slows down the initial accuracy gain during the very early phase of training, reflecting the expected delay when using outdated rollouts. However, this effect quickly diminishes: all curves steadily improve and converge to nearly the same accuracy. Notably, even under extreme staleness ($s = 256$), M2PO achieves performance on par with the on-policy case ($s = 0$), without signs of collapse or degradation. These results highlight that M2PO effectively preserves the learning signal

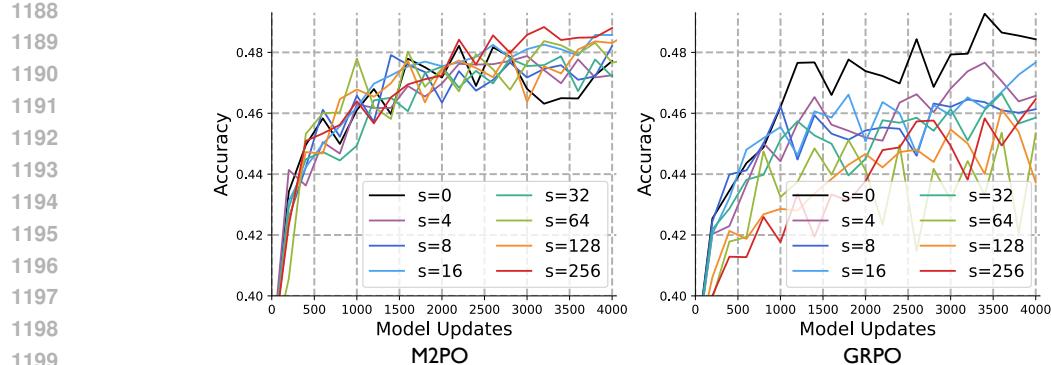


Figure 15: The performance of M2PO and GRPO under different staleness on Qwen2.5-Math-7B.

contained in stale data, enabling stable training and strong generalization across a wide range of staleness values.

G.4 TRAINING COLLAPSE ANALYSIS

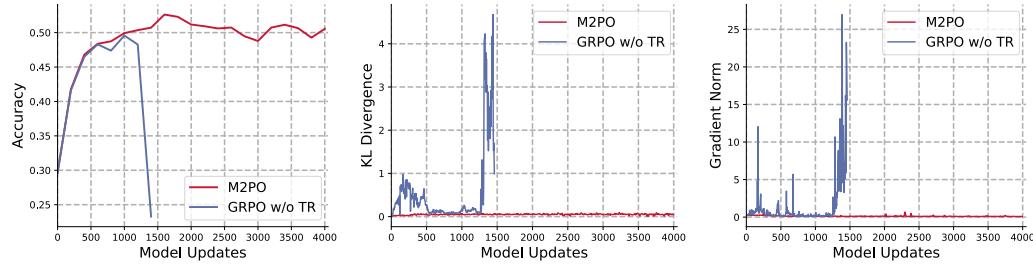


Figure 16: Qwen2.5-32B training dynamics comparison on M2PO and collapse training of GRPO without a trust region. **Left:** Test accuracy across eight math benchmarks. **Middle:** KL divergence between current policy and reference policy. **Right:** Gradient norm.

In our evaluation, we deem training to have collapsed when we observe a pronounced drop in test accuracy accompanied by abnormalities in stability metrics such as KL divergence and gradient norms. In practice, accuracy degradation is typically paired with sharp KL spikes and gradient-norm blow-ups. As shown in Figure 16, test accuracy drops at the same time that KL divergence and gradient norms exhibit consistent spikes.