
Reward-Wise Value Estimation for Multi-Reward Optimization in Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Value functions play a central role in critic-based reinforcement learning (RL)
2 algorithms such as proximal policy optimization (PPO) and in inference-time
3 scaling for large language models (LLMs). In complex LLM tasks, multiple
4 reward signals are often leveraged to guide RL training and inference. However,
5 existing value function estimation methods either aggregate all rewards into a
6 single scalar, limiting flexibility across reward components, or rely on separate
7 critic models for each reward, which is computationally inefficient. To address
8 these limitations, we propose an efficient and flexible value function estimation
9 method for multi-reward optimization in LLMs. Specifically, we introduce a multi-
10 head critic model architecture that jointly estimates reward-wise value functions
11 within a single model. This design enables fine-grained credit assignment across
12 reward components while maintaining computational efficiency. We further prove
13 that the proposed multi-head critic loss yields lower irreducible error compared
14 to the vanilla single-head critic loss. Experimental results on multi-turn agentic
15 search and mathematical reasoning tasks demonstrate that our reward-wise value
16 estimation approach consistently improves both policy optimization and critic
17 learning in PPO training.

18 1 Introduction

19 Reinforcement learning (RL) has become a standard paradigm for aligning large language models
20 (LLMs) on complex reasoning [1, 2] and agentic tasks [3, 4]. Existing RL algorithms for LLM training
21 can be broadly categorized into critic-free and critic-based methods. Critic-free algorithms, such as
22 Group Relative Policy Optimization (GRPO) [5] and REINFORCE Leave-One-Out (RLOO) [6, 7],
23 remove the value function and estimate advantages in a group-relative manner. While this avoids
24 the difficulty of critic training, the resulting advantage signals are typically derived from final
25 rewards and can be coarse-grained, leading to less precise credit assignment [8, 9, 10]. Critic-based
26 algorithms, such as proximal policy optimization (PPO) [11], leverage value functions to compute
27 generalized advantage estimation (GAE) [12], providing fine-grained advantage signals for policy
28 optimization [13]. Moreover, a well-trained critic can be used to guide inference-time scaling and
29 search [14, 15, 16, 17], further improving policy performance.

30 Recently, multiple reward signals have been widely used to guide RL training and inference in LLMs.
31 Beyond answer rewards, auxiliary rewards such as format rewards and process rewards are often
32 incorporated for reward shaping in search [18, 19], tool-calling agentic tasks [20, 21], preference
33 learning [22] and complex multi-agent intelligence [23]. However, these methods typically aggregate
34 all reward components into a single scalar, which limits flexibility across heterogeneous rewards
35 and may lead to unstable training or even training collapse [10, 24]. This naive aggregation can
36 lead to poor credit assignment, since the model loses access to reward-specific learning signals and

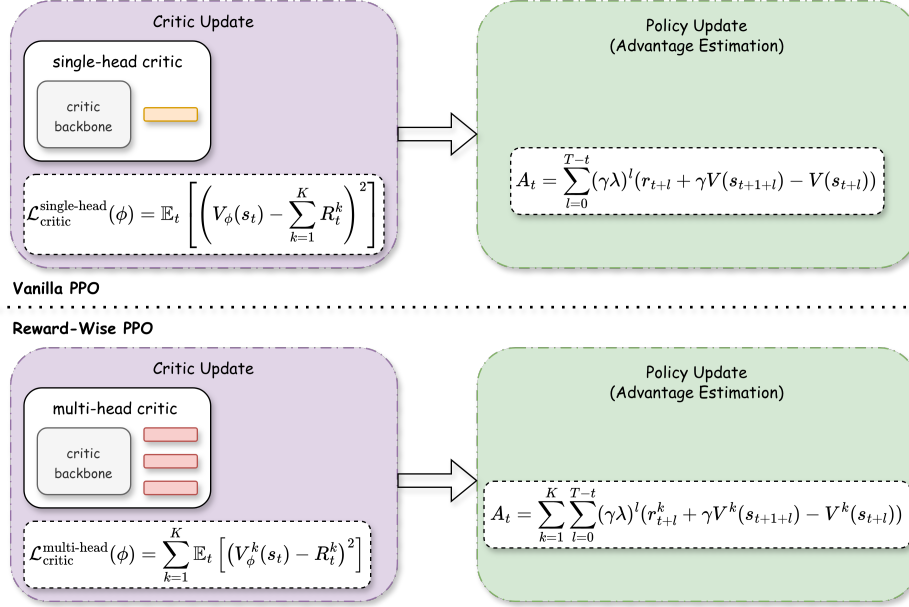


Figure 1: Comparison between vanilla PPO and reward-wise PPO with a multi-head critic. Vanilla PPO aggregates multiple reward components into a single scalar target and learns a single-head critic, while reward-wise PPO estimates reward-wise value functions with multiple value heads. The reward-wise advantages are computed separately for each reward component and then aggregated for policy optimization, enabling fine-grained value estimation across heterogeneous rewards with minimal additional critic overhead.

37 cannot identify which reward components contribute most during optimization. In addition, a critic
 38 trained on aggregated rewards can only estimate an overall value function, rather than providing
 39 reward-specific value signals for individual components. This limits its flexibility in inference-time
 40 scaling, where different reward components may need to be weighted or queried separately.

41 In this paper, we focus on critic-based RL algorithms for multi-reward settings. On the one hand,
 42 optimizing with aggregated rewards can lead to poor credit assignment and results in an inflexible
 43 critic that only provides overall value estimates, without reward-specific value signals for individual
 44 components. On the other hand, training a separate critic model for each reward component is
 45 computationally expensive and memory inefficient. Therefore, this raises the following question:

46 *How can we design critic-based RL algorithms that improve policy optimization while maintaining*
 47 *an efficient and flexible critic model for multi-reward settings?*

48 To address these limitations, we propose an efficient and flexible value function estimation method
 49 for multi-reward optimization in LLMs. Our contributions are summarized as follows:

- 50 • We introduce a multi-head critic architecture that jointly estimates reward-wise value func-
 51 tions within a single model. This design enables fine-grained value estimation across reward
 52 components while maintaining computational efficiency. We prove that the reward-wise
 53 critic loss yields lower irreducible error compared to the vanilla single-head critic loss.
- 54 • Building on the multi-head critic, we propose reward-wise PPO, which separately estimates
 55 GAE for each reward component and aggregates reward-wise advantages for policy optimiza-
 56 tion, enabling more fine-grained credit assignment across heterogeneous rewards. Figure 1
 57 compares vanilla PPO with a single-head critic and reward-wise PPO with a multi-head
 58 critic.
- 59 • Experimental results demonstrate that our reward-wise value estimation approach consis-
 60 tently improves both policy optimization and critic learning in PPO training across different
 61 model sizes and datasets on both multi-turn agentic search and math reasoning tasks.

62 2 Related Work

63 **RL Algorithms for LLMs.** For LLM training, existing RL algorithms can be categorized into critic-
64 free methods, such as GPRO [5] and RLOO [6, 7], and critic-based methods, such as PPO [11]. GRPO-
65 type methods estimate advantages from group-level rollout samples and avoid critic training, resulting
66 in lower memory and computational overhead. However, although they bypass the difficulty of
67 training a critic, their advantage signals are typically derived from final rewards and can lead to coarse-
68 grained trajectory-level advantage estimation, resulting in less precise credit assignment [8, 9, 10].
69 Recent studies further suggest that such coarse-grained advantage estimation can lead to unstable
70 training and even training collapse [10, 25, 26]. In contrast, PPO-type methods leverage value
71 functions to compute GAE [12], providing fine-grained credit assignment for policy optimization [13].
72 Moreover, a well-trained critic can be used to guide or evaluate inference-time scaling and search [14,
73 15, 16, 17], further improving performance.

74 **Multi-Reward RL.** Recently, RL training and inference-time scaling have moved beyond single-
75 reward optimization. Beyond answer correctness rewards, auxiliary rewards such as format rewards
76 and process rewards are often incorporated for reward shaping in search [18, 19], tool-calling agentic
77 tasks [20, 21], and complex multi-agent intelligence [23]. Another important application of multi-
78 reward RL is improving the efficiency of reasoning models while maintaining task performance,
79 primarily by introducing length-based rewards alongside outcome-based rewards [27, 28]. In addition,
80 preference learning methods can use different preference rubrics to adapt multiple reward modeling
81 to diverse tasks [22]. However, these methods typically aggregate all reward components into a single
82 scalar, which limits flexibility across heterogeneous rewards and may lead to unstable training or
83 even training collapse [10, 24]. The most closely related works are concurrent studies that propose
84 estimating advantages separately for each reward component in a GRPO-style manner [24, 29]. The
85 key difference is that these methods are critic-free, whereas our work focuses on critic-based RL.
86 Moreover, these GRPO-type methods do not learn critic models, and thus cannot directly provide
87 value functions for inference-time scaling or search.

88 3 Preliminaries

89 3.1 Modeling Language Generation as Token-Level MDP

90 Let $x \sim \mathcal{D}$ denote the input prompt from the dataset \mathcal{D} , and $y = [y_1, \dots, y_T]$ the generated response.
91 Text generation by a language model can be formulated as sampling from a policy in a Markov
92 decision process (MDP), defined as the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where \mathcal{S} denotes the state space, \mathcal{A}
93 the action space, P the transition dynamics, r the reward function, and γ the discount factor. At each
94 time step t , the state is $s_t = (x, y_{<t}) \in \mathcal{S}$, consisting of the prompt x and the previously generated
95 tokens, and the action $a_t = y_t \in \mathcal{A}$ corresponds to selecting the next token from the vocabulary.
96 The transition dynamics P are deterministic: given s_t and a_t , the next state is $s_{t+1} = (x, y_{\leq t})$. The
97 reward function $r_t = r(s_t, a_t)$ assigns a scalar value at each time step. Since it is a common practice
98 to use the discount factor $\gamma = 1$ in RL for LLM training, we omit γ in the subsequent notation
99 for simplicity, unless otherwise specified. The value function is defined as $V(s_t) = \mathbb{E} \left[\sum_{t'=t}^T r_{t'} \right]$,
100 representing the expected return from state s_t . RL aims to maximize the expected cumulative reward

$$J(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[\sum_{t=1}^T r_t \right]. \quad (1)$$

101 In this work, we consider a terminal reward provided only at the end of the sequence, denoted as
102 $\mathcal{R}(x, y)$, i.e., $r_t = 0$ for $t < T$ and $r_T = \mathcal{R}(x, y)$.

103 3.2 Proximal Policy Optimization

104 PPO [11] is a popular actor-critic RL algorithm commonly used for LLM training [13]. It optimizes
105 the policy by maximizing the clipped surrogate objective:

$$\mathcal{J}_{\text{PPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\sum_{t=1}^T \min(\rho_t A_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) A_t) \right], \quad (2)$$

106 where $\rho_t = \frac{\pi_\theta(y_t|y_{<t},x)}{\pi_{\theta_{\text{old}}}(y_t|y_{<t},x)}$ is the importance sampling ratio between the current policy π_θ and the old
 107 policy $\pi_{\theta_{\text{old}}}$, A_t is the advantage term, and ϵ is the clipping parameter. The PPO algorithm mainly
 108 consists of two components: advantage estimation and value estimation.

109 • **Advantage Estimation.** The advantage A_t measures the relative quality of action a_t at
 110 state s_t . In PPO, it is typically estimated using GAE [12]: $A_t = \sum_{l=0}^{T-t} (\gamma\lambda)^l \delta_{t+l}$, where
 111 $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ is the temporal-difference (TD) error, and $V(s_t)$ denotes the
 112 value of state s_t .

113 • **Value Estimation.** The value function, often referred to as the critic model, estimates the
 114 expected return for each state. The objective of the critic is to minimize the discrepancy
 115 between its predicted values and the target returns, typically via a mean squared error (MSE)
 116 loss:

$$\mathcal{L}_{\text{critic}}(\phi) = \mathbb{E}_t \left[(V_\phi(s_t) - R_t)^2 \right], \quad (3)$$

117 where $V_\phi(s_t)$ denotes the predicted value of state s_t parameterized by ϕ , and R_t denotes the
 118 target return, which can be estimated using Monte Carlo returns: $R_t = \sum_{t'=t}^T r_{t'}$.

119 3.3 Multi-Reward Modeling

120 In multi-reward modeling, the reward function in the MDP formulation outputs a K -dimensional
 121 vector, i.e., $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^K$, where r^k denotes the k -th reward component. The reward components
 122 are aggregated into a scalar reward, i.e., $r_t = \sum_{k=1}^K r_t^k$, for policy optimization in Eq. (1). When
 123 $K = 1$, this formulation reduces to the standard single-reward setting. In the terminal reward setting,
 124 each data pair (x, y) is associated with a set of K terminal reward components $\{\mathcal{R}^k(x, y)\}_{k=1}^K$. For
 125 multi-reward settings, learning an effective critic for policy optimization and inference-time scaling
 126 poses significant challenges. On the one hand, aggregating all reward components into a single scalar
 127 limits flexibility and obscures the contribution of individual rewards. On the other hand, learning
 128 separate critic models for each reward component incurs substantial computational and memory
 129 overhead.

130 4 Reward-Wise Value Estimation and Reward-Wise PPO

131 4.1 Reward-Wise Value Estimation

132 To address these limitations, we propose an efficient and flexible value function estimation method
 133 for multi-reward optimization in LLMs. Specifically, we introduce a multi-head critic architecture
 134 that jointly estimates reward-wise value functions within a single model. Compared to a traditional
 135 single-head critic, our multi-head critic shares the same backbone while incorporating multiple value
 136 heads, each learning a distinct reward component independently.

137 Given K Monte Carlo returns R_t^k for $k = 1, 2, \dots, K$, we define the reward-wise critic loss as

$$\mathcal{L}_{\text{critic}}^{\text{rw}}(\phi) = \sum_{k=1}^K \mathbb{E}_t \left[(V_\phi^k(s_t) - R_t^k)^2 \right], \quad (4)$$

138 where $V_\phi^k(s_t)$ denotes the predicted value of state s_t from the k -th value head corresponding to the
 139 k -th reward component. Compared to the vanilla critic loss

$$\mathcal{L}_{\text{critic}}(\phi) = \mathbb{E}_t \left[\left(V_\phi(s_t) - \sum_{k=1}^K R_t^k \right)^2 \right], \quad (5)$$

140 which estimates the aggregated return across all reward components, the proposed reward-wise critic
 141 loss decomposes value estimation into per-reward components, enabling more fine-grained modeling
 142 of heterogeneous rewards.

143 **Overhead Analysis.** We compare our multi-head critic architecture with the vanilla single-head critic
 144 and a multi-critic architecture in terms of parameter count and forward computational cost (FLOPs) in
 145 Table 1. For parameter count, all three critic architectures share the same backbone with $12lh^2 + vh$

Table 1: Parameter and forward FLOPs comparison of different critic architectures. Here, l denotes the number of Transformer layers, h the hidden dimension, v the vocabulary size, s the sequence length, b the batch size, and K the number of reward components. Detailed derivations of the parameter count and forward FLOPs can be found in Appendix A.

Architecture	Parameter Count	Forward FLOPs
single-head critic	$12lh^2 + vh + h$	$b[l(12sh^2 + 2s^2h) + sh]$
multi-critic	$K(12lh^2 + vh + h)$	$Kb[l(12sh^2 + 2s^2h) + sh]$
multi-head critic	$12lh^2 + vh + Kh$	$b[l(12sh^2 + 2s^2h) + Ksh]$

146 parameters, which dominates the overall parameter size. The difference arises from the output head:
 147 the single-head critic has a single value head with h parameters, the multi-critic architecture replicates
 148 the entire backbone K times, and the multi-head critic introduces only K lightweight value heads
 149 with Kh parameters. In particular, compared to the single-head critic, the multi-head critic introduces
 150 only an additional $(K - 1)h$ parameters, which corresponds to a negligible overhead relative to the
 151 backbone $12lh^2 + vh$. For forward FLOPs, all three critic architectures share the same backbone cost
 152 $l(12sh^2 + 2s^2h)$, which dominates the overall computation. The difference lies in the output head:
 153 the single-head critic incurs an additional sh cost, the multi-critic architecture requires K independent
 154 forward passes, and the multi-head critic introduces only Ksh additional cost from the value heads.
 155 Therefore, the multi-head critic achieves comparable parameter efficiency and computational cost to
 156 the single-head critic, while avoiding the K -fold increase in both parameters and FLOPs incurred by
 157 the multi-critic architecture.

158 We provide a theoretical analysis comparing the two critic losses in Eq. (4) and Eq. (5).

159 **Proposition 1.** *Let each reward-wise return be decomposed as*

$$R_t^k = \mu_t^k + \epsilon_t^k,$$

160 where $\mu_t^k = \mathbb{E}[R_t^k | s_t]$ and ϵ_t^k is zero-mean noise. Let $V_\phi^k(s_t)$ be the reward-wise value estimate and
 161 $V_\phi(s_t)$ be the vanilla value estimate for the aggregated return. Define the reward-wise prediction
 162 error as $e_t^k = V_\phi^k(s_t) - \mu_t^k$, and the vanilla prediction error as $e_t = V_\phi(s_t) - \sum_{k=1}^K \mu_t^k$. Then the
 163 expected reward-wise critic loss decomposes as

$$\mathbb{E}[\mathcal{L}_{\text{critic}}^{\text{rw}}] = \sum_{k=1}^K \mathbb{E}[(e_t^k)^2] + \sum_{k=1}^K \mathbb{E}[(\epsilon_t^k)^2],$$

164 whereas the expected vanilla critic loss decomposes as

$$\mathbb{E}[\mathcal{L}_{\text{critic}}] = \mathbb{E}[e_t^2] + \mathbb{E}\left[\left(\sum_{k=1}^K \epsilon_t^k\right)^2\right].$$

165 Moreover,

$$\mathbb{E}\left[\left(\sum_{k=1}^K \epsilon_t^k\right)^2\right] = \sum_{k=1}^K \mathbb{E}[(\epsilon_t^k)^2] + 2 \sum_{i < j} \mathbb{E}[\epsilon_t^i \epsilon_t^j].$$

166 Thus, the vanilla critic loss contains cross-component noise terms, while the reward-wise critic loss
 167 decomposes the regression problem across individual reward components. In particular, when the
 168 reward noise terms are positively correlated, i.e., $\mathbb{E}[\epsilon_t^i \epsilon_t^j] > 0$ for $i \neq j$, the cross-component terms
 169 further increase the noise level of the vanilla critic target.

170 The proof is provided in Appendix B. Here, the noise refers to the randomness in the Monte
 171 Carlo return targets induced by policy sampling, environment stochasticity, and reward estimation,
 172 formally defined as $\epsilon_t^k = R_t^k - \mathbb{E}[R_t^k | s_t]$. This corresponds to the irreducible noise in the critic
 173 regression target, which cannot be eliminated by improving the quality of value function estimation.
 174 When rewards are aggregated, these component-wise noise terms are also summed, introducing
 175 additional cross-component noise through covariance terms. In particular, when reward components
 176 are positively correlated, aggregation further amplifies the noise level in the vanilla critic target.

Algorithm 1: Reward-Wise PPO with Multi-Head Critic Learning

Input: Policy model π_θ , old policy $\pi_{\theta_{\text{old}}}$, multi-head critic $V_\phi = \{V_\phi^k\}_{k=1}^K$, reward components $\{r_t^k\}_{k=1}^K$

Output: Updated policy parameters θ and critic parameters ϕ

```
1 for each training iteration do
  // 1. Rollout stage
2   Collect trajectories using  $\pi_{\theta_{\text{old}}}$  and obtain reward components  $\{r_t^k\}_{k=1}^K$ ;
  // 2. Training stage
3   Compute reward-wise value estimates  $\{V_\phi^k(s_t)\}_{k=1}^K$  using the multi-head critic;
4   for each reward component  $k = 1, \dots, K$  do
5     Compute reward-wise advantages  $A_t^k$  using reward-wise GAE in Eq. (7);
6     Compute reward-wise returns  $R_t^k$ ;
7   Aggregate reward-wise advantages  $A_t = \sum_{k=1}^K A_t^k$  as in Eq. (6);
8   Update the policy using the PPO objective in Eq. (2) with the aggregated advantage  $A_t$ ;
9   Update the multi-head critic using the reward-wise critic loss in Eq. (4);
10  Set  $\theta_{\text{old}} \leftarrow \theta$ ;
```

177 In contrast, the reward-wise critic keeps the regression targets separated and therefore avoids such
178 aggregation-induced cross-component noise. As a result, it provides a cleaner supervision signal for
179 critic learning and mitigates interference across heterogeneous reward components.

180 4.2 Reward-Wise PPO

181 Building upon our reward-wise value estimation with a multi-head critic, we introduce reward-wise
182 PPO to further improve policy optimization. Specifically, we decompose the advantage function into
183 reward-wise components and aggregate them for policy updates.

184 We define the reward-wise advantage as

$$A_t = \sum_{k=1}^K A_t^k, \quad (6)$$

185 where each reward-wise advantage is estimated using GAE in a component-wise manner:

$$A_t^k = \sum_{l=0}^{T-t} (\gamma\lambda)^l \delta_{t+l}^k, \quad \delta_t^k = r_t^k + \gamma V_\phi^k(s_{t+1}) - V_\phi^k(s_t), \quad (7)$$

186 and $V_\phi^k(s_t)$ denotes the value estimate from the k -th value head corresponding to the k -th reward
187 component. The aggregated advantage A_t is then used in the PPO objective in Eq. (2).

188 Compared to vanilla PPO, which computes advantages from aggregated rewards, reward-wise PPO
189 estimates advantages separately for each reward component and aggregates them at the policy update
190 stage, enabling more fine-grained credit assignment across heterogeneous rewards. Please refer to the
191 algorithm pseudo-code of reward-wise PPO with multi-head critic learning in Algorithm 1 for details.

192 5 Experiments

193 We conduct experiments on multi-turn agentic search and mathematical reasoning tasks to demonstrate
194 the effectiveness of our method. We keep all hyperparameters and reward design the same for reward-
195 wise PPO (RW-PPO) and vanilla PPO. All experiments are conducted on 8×H100 GPUs. Due to
196 space limitations, experiments on mathematical reasoning tasks are provided in Appendix C.

197 5.1 Experimental Settings

198 **Models and Datasets.** We compare our proposed RW-PPO with vanilla PPO on multi-turn agentic
199 search tasks. We train Qwen2.5-3B-Instruct and Qwen2.5-7B-Instruct [30] on a mixed training

Table 2: Parameter and memory comparison of base models and critics. Memory is estimated in BF16 (2 bytes per parameter). Values in parentheses denote the difference compared to the 1-head critic. For Qwen2.5-3B, the LM head is tied with the input embedding, resulting in negligible parameter difference between the full model and the critic. In contrast, Qwen2.5-7B does not tie word embeddings, so removing the LM head leads to a significant reduction in parameters.

Model	Hidden Size	#Params	Memory (GB)
Qwen2.5-3B-Instruct (full)	2048	3,085,938,688	6.17
Critic (1-head)	2048	3,085,940,736	6.17
Critic (3-head)	2048	3,085,944,832 (+4,096)	6.17 (+8 KB)
Qwen2.5-7B-Instruct (full)	3584	7,615,616,512	15.23
Critic (1-head)	3584	7,070,622,720	14.14
Critic (3-head)	3584	7,070,629,888 (+7,168)	14.14 (+14 KB)

dataset constructed from NQ [31] and HotpotQA [32], and evaluate the trained models on the test sets of NQ [31], TriviaQA [33], PopQA [34], HotpotQA [32], 2Wiki [35], MuSiQue [36], and Bamboogle [37], following the Search-R1 setting [3]. We use E5 [38] as the retriever, and 2018 Wikipedia dump [31] as the corpus. We observe training instability and occasional collapse. For fair comparison, we report the maximum total reward achieved across all training steps.

Reward Design. We design three terminal reward components: (1) *Format correctness*. This reward ensures structural validity by verifying that the outputs in both the final turn and all intermediate turns comply with the specified format. (2) *Retrieval correctness*. This reward evaluates whether any accepted answer appears in at least one search result (from the <information> tag), using case-insensitive string matching. (3) *Answer correctness*. This reward evaluates whether the extracted answer (from the <answer> tag) exactly matches any accepted ground-truth answer after normalization (e.g., lowercasing and whitespace removal). We assign reward values as $\mathcal{R}_{\text{format}} \in \{0, 0.2\}$, $\mathcal{R}_{\text{ret}} \in \{0, 0.3\}$, and $\mathcal{R}_{\text{ans}} \in \{0, 1.0\}$.

PPO Training. We enable gradient checkpointing and use FSDP2 with CPU offloading. The learning rates of the policy and critic models are set to 1.0×10^{-6} and 1.0×10^{-5} , respectively. Training is performed for one epoch with one update epoch per batch. The total batch size is 512, with policy and critic mini-batch sizes of 256 and a micro-batch size of 4 per GPU. The maximum input length is set to 4,096 tokens, and the maximum generation length is set to 500 tokens. KL loss is disabled. We apply off-policy correction for token-level importance sampling, with the importance ratio clipped at 2.0. Specifically, we use RW-PPO with three value heads and vanilla PPO with a single value head for the critic.

Rollout Generation. We use vLLM with one inference engine per GPU, tensor parallel size 1, and a GPU memory utilization ratio of 0.5. We sample one response per prompt with a temperature of 1.0 and a top- p value of 1.0. We allow up to 4 search turns and retrieve the top 3 documents for each search query.

5.2 Overhead Analysis

As shown in Table 2, we compare the parameter and memory overhead of the multi-head critic with the single-head critic across different model scales. The multi-head critic introduces negligible overhead compared to the single-head critic. Specifically, adding multiple value heads increases the parameter count by only a few thousand parameters (e.g., +4,096 for 3B and +7,168 for 7B), corresponding to only several kilobytes of additional memory. This overhead is negligible relative to the total model size (on the order of billions of parameters and tens of gigabytes).

5.3 Performance Comparison

As shown in Table 3, RW-PPO consistently outperforms vanilla PPO across both model scales and most datasets. Compared to PPO, RW-PPO achieves higher total rewards on all datasets for the 3B model and on nearly all datasets for the 7B model, demonstrating the effectiveness of reward-wise optimization. These gains are primarily driven by improvements in answer reward, which show larger relative increases than other reward components. The improvements are particularly

Table 3: Performance comparison between RW-PPO and PPO on multi-turn agentic search tasks. We report Answer, Format, Retrieval, and Total rewards. Values in parentheses indicate relative improvements over PPO. Avg denotes the average across all datasets. † and * denote in-domain and out-of-domain datasets, respectively.

Dataset	Method	Answer (0–1)	Format (0–0.2)	Retrieval (0–0.3)	Total (0–1.5)
<i>Qwen2.5-3B-Instruct</i>					
NQ [†]	PPO	0.3825	0.1914	0.2171	0.7911
	RW-PPO	0.4371 (+14.3%)	0.1995	0.2171	0.8537 (+7.9%)
TriviaQA*	PPO	0.5717	0.1922	0.2229	0.9868
	RW-PPO	0.5985 (+4.7%)	0.1991	0.2246	1.0222 (+3.6%)
PopQA*	PPO	0.4046	0.1946	0.1976	0.7968
	RW-PPO	0.4281 (+5.8%)	0.1995	0.1972	0.8247 (+3.5%)
HotpotQA [†]	PPO	0.3618	0.1924	0.1828	0.7369
	RW-PPO	0.3904 (+7.9%)	0.1975	0.1873	0.7752 (+5.2%)
2Wiki*	PPO	0.3497	0.1757	0.1523	0.6777
	RW-PPO	0.3885 (+11.1%)	0.1889	0.1566	0.7340 (+8.3%)
MuSiQue*	PPO	0.1415	0.1764	0.1089	0.4268
	RW-PPO	0.1700 (+20.1%)	0.1820	0.1143	0.4663 (+9.3%)
Bamboogle*	PPO	0.3520	0.2000	0.1464	0.6984
	RW-PPO	0.4000 (+13.6%)	0.2000	0.1560	0.7560 (+8.2%)
Avg	PPO	0.3663	0.1889	0.1754	0.7306
	RW-PPO	0.4018 (+9.7%)	0.1952	0.1790	0.7760 (+6.2%)
<i>Qwen2.5-7B-Instruct</i>					
NQ [†]	PPO	0.4723	0.1987	0.2267	0.8977
	RW-PPO	0.4820 (+2.1%)	0.2000	0.2344	0.9164 (+2.1%)
TriviaQA*	PPO	0.6491	0.1979	0.2332	1.0802
	RW-PPO	0.6606 (+1.8%)	0.1999	0.2358	1.0962 (+1.5%)
PopQA*	PPO	0.4466	0.1992	0.1945	0.8403
	RW-PPO	0.4683 (+4.9%)	0.2000	0.2030	0.8713 (+3.7%)
HotpotQA [†]	PPO	0.4451	0.1965	0.2000	0.8416
	RW-PPO	0.4528 (+1.7%)	0.1998	0.2040	0.8566 (+1.8%)
2Wiki*	PPO	0.3985	0.1778	0.1654	0.7417
	RW-PPO	0.4369 (+9.6%)	0.1986	0.1606	0.7960 (+7.3%)
MuSiQue*	PPO	0.2156	0.1809	0.1393	0.5357
	RW-PPO	0.2069 (-4.0%)	0.1945	0.1343	0.5357 (+0.0%)
Bamboogle*	PPO	0.4480	0.2000	0.1608	0.8088
	RW-PPO	0.4640 (+3.6%)	0.1968	0.1536	0.8144 (+0.7%)
Avg	PPO	0.4393	0.1930	0.1886	0.8209
	RW-PPO	0.4531 (+3.1%)	0.1985	0.1894	0.8410 (+2.4%)

238 notable on out-of-domain datasets such as 2Wiki, MuSiQue, and Bamboogle, where RW-PPO yields
 239 gains of up to 9.3% in total reward and over 20% in answer reward for the 3B model, indicating
 240 better generalization ability. On in-domain datasets (NQ and HotpotQA), RW-PPO also brings
 241 consistent improvements, though the gains are relatively smaller, suggesting that the benefits are
 242 more pronounced under distribution shift. Overall, RW-PPO improves the average total reward by
 243 6.2% on the 3B model and 2.4% on the 7B model, while achieving larger gains in answer reward
 244 (+9.7% and +3.1%, respectively), showing stronger advantages on smaller models.

245 As shown in Figure 2, RW-PPO achieves more stable reward trajectories than vanilla PPO on multi-
 246 turn agentic search tasks. For the 3B model, PPO collapses in answer, format, and retrieval rewards,
 247 while RW-PPO keeps these reward components stable. The total return also remains higher under
 248 RW-PPO, suggesting that the multi-head critic learns better value estimates for advantage estimation.
 249 For the 7B model, training is more stable for both methods, but RW-PPO still yields higher answer
 250 reward, retrieval reward, and total return.

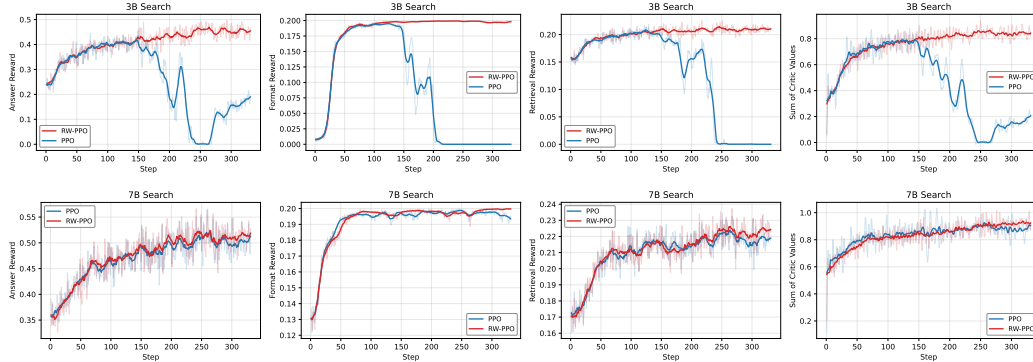


Figure 2: Training dynamics on multi-turn agentic search tasks. The top row shows results for Qwen2.5-3B-Instruct, and the bottom row shows results for Qwen2.5-7B-Instruct. From left to right, we report answer reward, format reward, retrieval reward, and total critic values during training. Compared with vanilla PPO, RW-PPO yields more stable reward trajectories and consistently improves the final total reward, especially for the 3B model.

Table 4: Answer reward comparison across methods. We compare GRPO and PPO with a single answer reward in Search-R1, PPO with multiple rewards, and RW-PPO with multiple rewards on multi-turn agentic search benchmarks. [†] and * denote in-domain and out-of-domain datasets, respectively.

Method	NQ [†]	TriviaQA*	PopQA*	HotpotQA [†]	2Wiki*	MuSiQue*	Bamboogle*	Avg
<i>Qwen2.5-3B-Instruct</i>								
GRPO (single-reward)	0.397	0.565	0.391	0.331	0.310	0.124	0.232	0.336
PPO (single-reward)	0.341	0.545	0.378	0.324	0.319	0.103	0.264	0.325
PPO (multi-reward)	0.383	0.572	0.405	0.362	0.350	0.142	0.352	0.366
RW-PPO (multi-reward)	0.437	0.599	0.428	0.390	0.389	0.170	0.400	0.402
<i>Qwen2.5-7B-Instruct</i>								
GRPO (single-reward)	0.429	0.623	0.427	0.386	0.346	0.162	0.400	0.396
PPO (single-reward)	0.393	0.610	0.397	0.370	0.414	0.146	0.368	0.385
PPO (multi-reward)	0.472	0.649	0.447	0.445	0.399	0.216	0.448	0.439
RW-PPO (multi-reward)	0.482	0.661	0.468	0.453	0.437	0.207	0.464	0.453

251 5.4 Benchmark with Other RL Methods

252 We further benchmark against GRPO and PPO with a single answer reward under the Search-R1
 253 setting [18, 3]. As shown in Table 4, using multiple rewards consistently improves over single-
 254 reward training, showing the benefit of reward shaping. RW-PPO further achieves the best average
 255 performance on both 3B and 7B models, demonstrating that reward-wise value estimation can better
 256 leverage heterogeneous reward signals.

257 6 Conclusion

258 In this paper, we propose reward-wise PPO with a multi-head critic for multi-reward RL in LLMs,
 259 enabling reward-wise value estimation and advantage computation while sharing a single efficient
 260 critic backbone. Our analysis shows that reward-wise critic learning avoids aggregation-induced
 261 cross-component noise, and experiments on multi-turn agentic search and mathematical reasoning
 262 tasks demonstrate consistent improvements in both answer reward and total reward over vanilla
 263 PPO. These results highlight reward-wise value estimation as an efficient and flexible approach for
 264 critic-based multi-reward optimization. One limitation of our work is that our experiments consider
 265 a relatively small number of reward components. An interesting future direction is to extend our
 266 method to large-scale multi-reward RL scenarios with hundreds of reward components.

267 **References**

- 268 [1] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu,
269 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
270 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 271 [2] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He.
272 Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the
273 wild. *arXiv preprint arXiv:2503.18892*, 2025.
- 274 [3] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Za-
275 mani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with
276 reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
- 277 [4] Zhenghai Xue, Longtao Zheng, Qian Liu, Yingru Li, Xiaosen Zheng, Zejun Ma, and Bo An.
278 Simpletir: End-to-end reinforcement learning for multi-turn tool-integrated reasoning. *arXiv*
279 *preprint arXiv:2509.02479*, 2025.
- 280 [5] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
281 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
282 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 283 [6] Wouter Kool, Herke van Hoof, and Max Welling. Buy 4 reinforce samples, get a baseline for
284 free! 2019.
- 285 [7] Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier
286 Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce-style optimization
287 for learning from human feedback in llms. In *Proceedings of the 62nd Annual Meeting of the*
288 *Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12248–12267, 2024.
- 289 [8] Yiran Guo, Lijie Xu, Jie Liu, Dan Ye, and Shuang Qiu. Segment policy optimization: Ef-
290 fective segment-level credit assignment in rl for large language models. *arXiv preprint*
291 *arXiv:2505.23564*, 2025.
- 292 [9] Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. Group-in-group policy optimization for
293 llm agent training. *arXiv preprint arXiv:2505.10978*, 2025.
- 294 [10] Quan Wei, Siliang Zeng, Chenliang Li, William Brown, Oana Frunza, Wei Deng, Anderson
295 Schneider, Yuriy Nevmyvaka, Yang Katie Zhao, Alfredo Garcia, and Mingyi Hong. Reinforcing
296 multi-turn reasoning in llm agents via turn-level reward design, 2025.
- 297 [11] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
298 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 299 [12] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-
300 dimensional continuous control using generalized advantage estimation. *arXiv preprint*
301 *arXiv:1506.02438*, 2015.
- 302 [13] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
303 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
304 follow instructions with human feedback. *Advances in neural information processing systems*,
305 35:27730–27744, 2022.
- 306 [14] Zhixuan Liu, Zhanhui Zhou, Yuanfu Wang, Chao Yang, and Yu Qiao. Inference-time language
307 model alignment via integrated value guidance. In *Findings of the Association for Computational*
308 *Linguistics: EMNLP 2024*, pages 4181–4195, 2024.
- 309 [15] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
310 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-*
311 *vances in neural information processing systems*, 36:11809–11822, 2023.
- 312 [16] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee,
313 Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The*
314 *twelfth international conference on learning representations*, 2023.

- 315 [17] Amrith Setlur, Chirag Nagpal, Adam Fisch, Xinyang Geng, Jacob Eisenstein, Rishabh Agarwal,
316 Alekh Agarwal, Jonathan Berant, and Aviral Kumar. Rewarding progress: Scaling automated
317 process verifiers for llm reasoning. *arXiv preprint arXiv:2410.08146*, 2024.
- 318 [18] Bowen Jin, Jinsung Yoon, Priyanka Kargupta, Sercan O Arik, and Jiawei Han. An empirical
319 study on reinforcement learning for reasoning-search interleaved llm agents. *arXiv preprint*
320 *arXiv:2505.15117*, 2025.
- 321 [19] Shu Zhao, Tan Yu, Anbang Xu, Japinder Singh, Aaditya Shukla, and Rama Akkiraju. Par-
322 allelsearch: Train your llms to decompose query and search sub-queries in parallel with
323 reinforcement learning. *arXiv preprint arXiv:2508.09303*, 2025.
- 324 [20] Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *arXiv preprint*
325 *arXiv:2503.23383*, 2025.
- 326 [21] Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan
327 Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*,
328 2025.
- 329 [22] Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar Nath, Yunzhong He, Bing Liu, and Sean
330 Hendryx. Rubrics as rewards: Reinforcement learning beyond verifiable domains. *arXiv*
331 *preprint arXiv:2507.17746*, 2025.
- 332 [23] Kimi Team, Tongtong Bai, Yifan Bai, Yiping Bao, S. H. Cai, Yuan Cao, Y. Charles, H. S. Che,
333 Cheng Chen, Guanduo Chen, Huarong Chen, Jia Chen, Jiahao Chen, Jianlong Chen, Jun Chen,
334 Kefan Chen, Liang Chen, Ruijue Chen, Xinhao Chen, Yanru Chen, Yanxu Chen, Yicun Chen,
335 Yimin Chen, Yingjiang Chen, Yuankun Chen, Yujie Chen, Yutian Chen, Zhirong Chen, Ziwei
336 Chen, Dazhi Cheng, Minghan Chu, Jialei Cui, Jiaqi Deng, Muxi Diao, Hao Ding, Mengfan
337 Dong, Mengnan Dong, Yuxin Dong, Yuhao Dong, Angang Du, Chenzhuang Du, Dikang Du,
338 Lingxiao Du, Yulun Du, Yu Fan, Shengjun Fang, Qiulin Feng, Yichen Feng, Garimugai Fu,
339 Kelin Fu, Hongcheng Gao, Tong Gao, Yuyao Ge, Shangyi Geng, Chengyang Gong, Xiaochen
340 Gong, Zhuoma Gongque, Qizheng Gu, Xinran Gu, Yicheng Gu, Longyu Guan, Yuanying Guo,
341 Xiaoru Hao, Weiran He, Wenyang He, Yunjia He, Chao Hong, Hao Hu, Jiayi Hu, Yangyang Hu,
342 Zhenxing Hu, Ke Huang, Ruiyuan Huang, Weixiao Huang, Zhiqi Huang, Tao Jiang, Zhejun
343 Jiang, Xinyi Jin, Yu Jing, Guokun Lai, Aidi Li, C. Li, Cheng Li, Fang Li, Guanghe Li, Guanyu
344 Li, Haitao Li, Haoyang Li, Jia Li, Jingwei Li, Junxiong Li, Lincan Li, Mo Li, Weihong Li,
345 Wentao Li, Xinhao Li, Xinhao Li, Yang Li, Yanhao Li, Yiwei Li, Yuxiao Li, Zhaowei Li,
346 Zhengming Li, Weilong Liao, Jiawei Lin, Xiaohan Lin, Zhishan Lin, Zichao Lin, Cheng Liu,
347 Chenyu Liu, Hongzhang Liu, Liang Liu, Shaowei Liu, Shudong Liu, Shuran Liu, Tianwei Liu,
348 Tianyu Liu, Weizhou Liu, Xiangyan Liu, Yangyang Liu, Yanming Liu, Yibo Liu, Yuanxin Liu,
349 Yue Liu, Zhengying Liu, Zhongnuo Liu, Enzhe Lu, Haoyu Lu, Zhiyuan Lu, Junyu Luo, Tongxu
350 Luo, Yashuo Luo, Long Ma, Yingwei Ma, Shaoguang Mao, Yuan Mei, Xin Men, Fanqing Meng,
351 Zhiyong Meng, Yibo Miao, Mingqing Ni, Kun Ouyang, Siyuan Pan, Bo Pang, Yuchao Qian,
352 Ruoyu Qin, Zeyu Qin, Jiezhong Qiu, Bowen Qu, Zeyu Shang, Youbo Shao, Tianxiao Shen,
353 Zhennan Shen, Juanfeng Shi, Lidong Shi, Shengyuan Shi, Feifan Song, Pengwei Song, Tianhui
354 Song, Xiaoxi Song, Hongjin Su, Jianlin Su, Zhaochen Su, Lin Sui, Jinsong Sun, Junyao Sun,
355 Tongyu Sun, Flood Sung, Yunpeng Tai, Chuning Tang, Heyi Tang, Xiaojuan Tang, Zhengyang
356 Tang, Jiawen Tao, Shiyuan Teng, Chaoran Tian, Pengfei Tian, Ao Wang, Bowen Wang, Chensi
357 Wang, Chuang Wang, Congcong Wang, Dingkun Wang, Dinglu Wang, Dongliang Wang, Feng
358 Wang, Hailong Wang, Haiming Wang, Hengzhi Wang, Huaqing Wang, Hui Wang, Jiahao Wang,
359 Jinhong Wang, Jiuzheng Wang, Kaixin Wang, Linian Wang, Qibin Wang, Shengjie Wang, Shuyi
360 Wang, Si Wang, Wei Wang, Xiaochen Wang, Xinyuan Wang, Yao Wang, Yejie Wang, Yipu
361 Wang, Yiqin Wang, Yucheng Wang, Yuzhi Wang, Zhaoji Wang, Zhaowei Wang, Zhengtao
362 Wang, Zhexu Wang, Zihan Wang, Zizhe Wang, Chu Wei, Ming Wei, Chuan Wen, Zichen Wen,
363 Chengjie Wu, Haoning Wu, Junyan Wu, Rucong Wu, Wenhao Wu, Yuefeng Wu, Yuhao Wu,
364 Yuxin Wu, Zijian Wu, Chenjun Xiao, Jin Xie, Xiaotong Xie, Yuchong Xie, Yifei Xin, Bowei
365 Xing, Boyu Xu, Jianfan Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu,
366 Xinbo Xu, Xinran Xu, Yangchuan Xu, Yichang Xu, Yuemeng Xu, Zelai Xu, Ziyao Xu, Junjie
367 Yan, Yuzi Yan, Guangyao Yang, Hao Yang, Junwei Yang, Kai Yang, Ningyuan Yang, Ruihan
368 Yang, Xiaofei Yang, Xinlong Yang, Ying Yang, Yi Yang, Yi Yang, Zhen Yang, Zhilin Yang,
369 Zonghan Yang, Haotian Yao, Dan Ye, Wenjie Ye, Zhuorui Ye, Bohong Yin, Chengzhen Yu,

- 370 Longhui Yu, Tao Yu, Tianxiang Yu, Enming Yuan, Mengjie Yuan, Xiaokun Yuan, Yang Yue,
371 Weihao Zeng, Dunyuan Zha, Haobing Zhan, Dehao Zhang, Hao Zhang, Jin Zhang, Puqi Zhang,
372 Qiao Zhang, Rui Zhang, Xiaobin Zhang, Y. Zhang, Yadong Zhang, Yangkun Zhang, Yichi
373 Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yushun Zhang, Yutao Zhang, Yutong Zhang,
374 Zheng Zhang, Chenguang Zhao, Feifan Zhao, Jinxiang Zhao, Shuai Zhao, Xiangyu Zhao, Yikai
375 Zhao, Zijia Zhao, Huabin Zheng, Ruihan Zheng, Shaojie Zheng, Tengyang Zheng, Junfeng
376 Zhong, Longguang Zhong, Weiming Zhong, M. Zhou, Runjie Zhou, Xinyu Zhou, Zaida Zhou,
377 Jinguo Zhu, Liya Zhu, Xinhao Zhu, Yuxuan Zhu, Zhen Zhu, Jingze Zhuang, Weiyu Zhuang,
378 Ying Zou, and Xinxing Zu. Kimi k2.5: Visual agentic intelligence, 2026.
- 379 [24] Shih-Yang Liu, Xin Dong, Ximing Lu, Shizhe Diao, Peter Belcak, Mingjie Liu, Min-Hung
380 Chen, Hongxu Yin, Yu-Chiang Frank Wang, Kwang-Ting Cheng, et al. Gdpo: Group reward-
381 decoupled normalization policy optimization for multi-reward rl optimization. *arXiv preprint*
382 *arXiv:2601.05242*, 2026.
- 383 [25] Yingru Li, Jiawei Xu, Ziniu Li, Jiakai Liu, Wei Liu, Yuxuan Tong, Longtao Zheng, Zhenghai
384 Xue, Yaxiang Zhang, Tianle Cai, et al. The optimal token baseline: Variance reduction for
385 long-horizon llm-rl. *arXiv preprint arXiv:2602.07078*, 2026.
- 386 [26] Zihan Wang, Kangrui Wang, Qineng Wang, Pingyue Zhang, Linjie Li, Zhengyuan Yang, Xing
387 Jin, Kefan Yu, Minh Nhat Nguyen, Licheng Liu, et al. Ragen: Understanding self-evolution in
388 llm agents via multi-turn reinforcement learning. *arXiv preprint arXiv:2504.20073*, 2025.
- 389 [27] Wei Liu, Ruochen Zhou, Yiyun Deng, Yuzhen Huang, Junteng Liu, Yuntian Deng, Yizhe Zhang,
390 and Junxian He. Learn to reason efficiently with adaptive length-based reward shaping. *arXiv*
391 *preprint arXiv:2505.15612*, 2025.
- 392 [28] Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with
393 reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.
- 394 [29] Jie Zhu, Yiyang Su, and Xiaoming Liu. Can textual reasoning improve the performance of
395 mllms on fine-grained visual classification? *arXiv preprint arXiv:2601.06993*, 2026.
- 396 [30] Qwen Team. Qwen2.5: A party of foundation models, September 2024.
- 397 [31] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov,
398 Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering.
399 In *EMNLP (1)*, pages 6769–6781, 2020.
- 400 [32] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhut-
401 dinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop
402 question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- 403 [33] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale
404 distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and
405 Min-Yen Kan, editors, *Proceedings of the 55th Annual Meeting of the Association for Computa-*
406 *tional Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada, July 2017.
407 Association for Computational Linguistics.
- 408 [34] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Ha-
409 jishirzi. When not to trust language models: Investigating effectiveness of parametric and
410 non-parametric memories. *arXiv preprint arXiv:2212.10511*, 2022.
- 411 [35] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing
412 a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint*
413 *arXiv:2011.01060*, 2020.
- 414 [36] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique:
415 Multihop questions via single-hop question composition. *Transactions of the Association for*
416 *Computational Linguistics*, 10:539–554, 2022.
- 417 [37] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis.
418 Measuring and narrowing the compositionality gap in language models. In *Findings of the*
419 *Association for Computational Linguistics: EMNLP 2023*, pages 5687–5711, 2023.

- 420 [38] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan
421 Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training.
422 *arXiv preprint arXiv:2212.03533*, 2022.
- 423 [39] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement
424 learning, 2025.
- 425 [40] Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai,
426 Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement
427 learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.

428 **A Detailed Derivations of Parameter Count and Forward FLOPs**

429 In this section, we provide detailed derivations of the parameter count and forward computational
 430 cost (FLOPs) for a Transformer-based critic model.

431 **A.1 Parameter Count**

432 We consider a Transformer-based critic with hidden dimension h , vocabulary size v , and l Transformer
 433 layers. The total parameter count consists of three parts: input embedding, Transformer backbone,
 434 and output head.

Table 5: Parameter count of each component in a Transformer-based critic. Here, h denotes the hidden dimension, v the vocabulary size, and l the number of layers.

Component	Parameter Count
Input Embedding	vh
Q, K, V projections	$3h^2$
Output projection	h^2
Feed-Forward Network	$8h^2$
Per-layer total	$12h^2$
Backbone (l layers)	$12lh^2$
Output Head	h
Total	$12lh^2 + vh + h$

435 The attention module consists of four projection matrices (query, key, value, and output), contributing
 436 $4h^2$ parameters per layer. The feed-forward network (FFN) expands the hidden dimension from h to
 437 $4h$ and projects it back, contributing $8h^2$ parameters. Therefore, each layer has $12h^2$ parameters, and
 438 the total backbone parameter count is $12lh^2$. The embedding layer contributes vh parameters, and
 439 the value head contributes h parameters.

440 **A.2 Forward FLOPs**

441 We now derive the forward computational cost (FLOPs). Let b denote the batch size and s the
 442 sequence length.

Table 6: Forward FLOPs of each component in a Transformer-based critic. Here, b denotes the batch size, s the sequence length, h the hidden dimension, and l the number of layers.

Component	Forward FLOPs
Input Embedding	negligible
Q, K, V projections	$3bsh^2$
Output projection	bsh^2
Attention computation	$2bs^2h$
Feed-Forward Network	$8bsh^2$
Per-layer total	$12bsh^2 + 2bs^2h$
Backbone (l layers)	$bl(12sh^2 + 2s^2h)$
Output Head	bsh
Total	$b [l(12sh^2 + 2s^2h) + sh]$

443 Although the embedding layer can be written as a matrix multiplication with one-hot inputs, it
 444 is implemented as a lookup operation in practice and thus incurs negligible FLOPs. For each
 445 Transformer layer, the main computational cost comes from linear projections, attention computation,
 446 and the feed-forward network, resulting in $12bsh^2 + 2bs^2h$ FLOPs per layer. Therefore, the total
 447 backbone cost is $bl(12sh^2 + 2s^2h)$. The output head introduces an additional bsh FLOPs.

448 **A.3 Discussion**

449 From Table 5 and Table 6, we observe that both the parameter count and computational cost are
 450 dominated by the Transformer backbone, particularly the $O(lh^2)$ parameters and $O(blsh^2)$ FLOPs.
 451 In contrast, the output head contributes negligible overhead. This observation forms the basis for the
 452 efficiency of the proposed multi-head critic design.

453 **B Proof of Proposition 1**

454 *Proof.* For the reward-wise critic loss, we have

$$\mathcal{L}_{\text{critic}}^{\text{rw}} = \sum_{k=1}^K (V_{\phi}^k(s_t) - R_t^k)^2.$$

455 Substituting $R_t^k = \mu_t^k + \epsilon_t^k$ yields

$$\mathcal{L}_{\text{critic}}^{\text{rw}} = \sum_{k=1}^K (V_{\phi}^k(s_t) - \mu_t^k - \epsilon_t^k)^2 = \sum_{k=1}^K (e_t^k - \epsilon_t^k)^2.$$

456 Taking expectation and using $\mathbb{E}[\epsilon_t^k | s_t] = 0$, we obtain

$$\mathbb{E}[\mathcal{L}_{\text{critic}}^{\text{rw}}] = \sum_{k=1}^K \mathbb{E}[(e_t^k)^2] + \sum_{k=1}^K \mathbb{E}[(\epsilon_t^k)^2].$$

457 For the vanilla critic loss, the aggregated return is

$$R_t = \sum_{k=1}^K R_t^k = \sum_{k=1}^K \mu_t^k + \sum_{k=1}^K \epsilon_t^k.$$

458 Thus,

$$\mathcal{L}_{\text{critic}} = \left(V_{\phi}(s_t) - \sum_{k=1}^K R_t^k \right)^2 = \left(e_t - \sum_{k=1}^K \epsilon_t^k \right)^2.$$

459 Taking expectation and using the zero-mean noise assumption, we obtain

$$\mathbb{E}[\mathcal{L}_{\text{critic}}] = \mathbb{E}[e_t^2] + \mathbb{E} \left[\left(\sum_{k=1}^K \epsilon_t^k \right)^2 \right].$$

460 Expanding the squared sum yields

$$\mathbb{E} \left[\left(\sum_{k=1}^K \epsilon_t^k \right)^2 \right] = \sum_{k=1}^K \mathbb{E}[(\epsilon_t^k)^2] + 2 \sum_{i < j} \mathbb{E}[\epsilon_t^i \epsilon_t^j],$$

461 which completes the proof. □

462 **C Experiments on Math Reasoning**

463 **C.1 Experimental Settings**

464 **Experimental Setup.** We compare our proposed RW-PPO with vanilla PPO on mathematical
 465 reasoning tasks. We train DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-7B [39]
 466 on the DAPO-Math-17k [40] dataset, and evaluate them on AIME 2024, AIME 2025, AIME 2026,
 467 and AMC datasets. The maximum sequence length is set to 8k tokens.

468 **Reward Design.** We design three terminal reward components: (1) *Format correctness*: the model is
 469 required to produce outputs with the correct structure, i.e., using XML tags to separate the reasoning
 470 (<think>) and the final answer (<answer>). (2) *Integer correctness*: whether the predicted answer

471 is a valid integer. (3) *Answer correctness*: whether the final answer matches the ground-truth answer.
472 We set the reward values as $\mathcal{R}_{\text{format}} \in \{0, 1.0\}$, $\mathcal{R}_{\text{int}} \in \{0, 1.0\}$, and $\mathcal{R}_{\text{ans}} \in \{0, 1.0\}$.

473 **PPO Training.** We enable gradient checkpointing and use FSDP with parameter and optimizer
474 offloading. The learning rates of the policy and critic models are set to 1.0×10^{-6} and 1.0×10^{-5} ,
475 respectively. Training is performed for one epoch. The total batch size is 1,024, with a mini-batch
476 size of 256 and a micro-batch size of 4 per GPU. The maximum prompt length is set to 2,048 tokens,
477 and the maximum response length is set to 8,192 tokens. KL loss is disabled. We apply token-level
478 importance sampling for rollout correction, with the importance ratio clipped at 2.0. Specifically, we
479 use RW-PPO with three value heads and vanilla PPO with a single value head for the critic.

480 **Rollout Generation.** We use vLLM for rollout generation with tensor parallel size 1 and a GPU
481 memory utilization ratio of 0.5. During training, we sample one response per prompt with a tempera-
482 ture of 1.0 and a top- p value of 1.0. During validation, we sample 16 responses per prompt with a
483 temperature of 0.7 and a top- p value of 0.95.

484 C.2 Experiment Results

485 As shown in Table 7, RW-PPO consistently outperforms vanilla PPO across all datasets and model
486 scales. RW-PPO achieves higher total rewards on every benchmark for both the 1.5B and 7B models,
487 demonstrating the effectiveness of reward-wise optimization in mathematical reasoning tasks. These
488 improvements are primarily driven by gains in answer reward, which show significantly larger relative
489 increases compared to format and integer rewards that are already near saturation. The improvements
490 are especially pronounced on more challenging benchmarks such as AIME 2025 and AIME 2026,
491 where RW-PPO yields substantial gains in answer reward (e.g., up to 29.0% for the 1.5B model),
492 leading to consistent improvements in total reward. In contrast, format and integer rewards remain
493 close to their upper bounds for both methods, indicating that RW-PPO mainly improves reasoning
494 accuracy rather than surface-level correctness. Overall, RW-PPO improves the average total reward
495 by 2.3% on the 1.5B model and 1.6% on the 7B model, while achieving larger gains in answer reward
496 (+11.3% and +5.5%, respectively), showing stronger benefits on smaller models.

Table 7: Performance comparison between RW-PPO and PPO on mathematical reasoning benchmarks. We report Answer, Format, Integer, and Total rewards averaged over 16 samples. Values in parentheses indicate relative improvements over PPO.

Dataset	Method	Answer (0–1)	Format (0–1)	Integer (0–1)	Total (0–3)
<i>DeepSeek-R1-Distill-Qwen-1.5B</i>					
AIME 2024	PPO	0.1396	0.9563	0.9563	2.0521
	RW-PPO	0.1521 (+8.9%)	0.9667	0.9667	2.0854 (+1.6%)
AIME 2025	PPO	0.1583	0.9854	0.9833	2.1271
	RW-PPO	0.2042 (+29.0%)	0.9979	0.9979	2.2000 (+3.4%)
AIME 2026	PPO	0.1500	0.9729	0.9729	2.0958
	RW-PPO	0.1667 (+11.1%)	0.9917	0.9896	2.1479 (+2.5%)
AMC	PPO	0.6516	0.9922	0.9906	2.6344
	RW-PPO	0.7000 (+7.4%)	0.9922	0.9922	2.6844 (+1.9%)
Avg	PPO	0.2749	0.9767	0.9758	2.2274
	RW-PPO	0.3058 (+11.3%)	0.9871	0.9866	2.2794 (+2.3%)
<i>DeepSeek-R1-Distill-Qwen-7B</i>					
AIME 2024	PPO	0.2313	0.9875	0.9875	2.2062
	RW-PPO	0.2500 (+8.1%)	0.9958	0.9958	2.2417 (+1.6%)
AIME 2025	PPO	0.2792	0.9938	0.9938	2.2667
	RW-PPO	0.3083 (+10.4%)	1.0000	1.0000	2.3083 (+1.8%)
AIME 2026	PPO	0.3229	0.9750	0.9750	2.2729
	RW-PPO	0.3583 (+11.0%)	0.9917	0.9917	2.3417 (+3.0%)
AMC	PPO	0.8281	0.9984	0.9984	2.8250
	RW-PPO	0.8359 (+0.9%)	1.0000	0.9984	2.8344 (+0.3%)
Avg	PPO	0.4154	0.9887	0.9887	2.3927
	RW-PPO	0.4381 (+5.5%)	0.9969	0.9965	2.4315 (+1.6%)

497 **NeurIPS Paper Checklist**

498 **1. Claims**

499 Question: Do the main claims made in the abstract and introduction accurately reflect the
500 paper’s contributions and scope?

501 Answer: [\[Yes\]](#)

502 Justification: The abstract and introduction clearly state the claims made

503 Guidelines:

- 504 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims
505 made in the paper.
- 506 • The abstract and/or introduction should clearly state the claims made, including the
507 contributions made in the paper and important assumptions and limitations. A [\[No\]](#) or
508 [\[N/A\]](#) answer to this question will not be perceived well by the reviewers.
- 509 • The claims made should match theoretical and experimental results, and reflect how
510 much the results can be expected to generalize to other settings.
- 511 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
512 are not attained by the paper.

513 **2. Limitations**

514 Question: Does the paper discuss the limitations of the work performed by the authors?

515 Answer: [\[Yes\]](#)

516 Justification: The paper has limitations discussed in the conclusion part in Section 6

517 Guidelines:

- 518 • The answer [\[N/A\]](#) means that the paper has no limitation while the answer [\[No\]](#) means
519 that the paper has limitations, but those are not discussed in the paper.
- 520 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 521 • The paper should point out any strong assumptions and how robust the results are to
522 violations of these assumptions (e.g., independence assumptions, noiseless settings,
523 model well-specification, asymptotic approximations only holding locally). The authors
524 should reflect on how these assumptions might be violated in practice and what the
525 implications would be.
- 526 • The authors should reflect on the scope of the claims made, e.g., if the approach was
527 only tested on a few datasets or with a few runs. In general, empirical results often
528 depend on implicit assumptions, which should be articulated.
- 529 • The authors should reflect on the factors that influence the performance of the approach.
530 For example, a facial recognition algorithm may perform poorly when image resolution
531 is low or images are taken in low lighting. Or a speech-to-text system might not be
532 used reliably to provide closed captions for online lectures because it fails to handle
533 technical jargon.
- 534 • The authors should discuss the computational efficiency of the proposed algorithms
535 and how they scale with dataset size.
- 536 • If applicable, the authors should discuss possible limitations of their approach to
537 address problems of privacy and fairness.
- 538 • While the authors might fear that complete honesty about limitations might be used by
539 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
540 limitations that aren’t acknowledged in the paper. The authors should use their best
541 judgment and recognize that individual actions in favor of transparency play an impor-
542 tant role in developing norms that preserve the integrity of the community. Reviewers
543 will be specifically instructed to not penalize honesty concerning limitations.

544 **3. Theory assumptions and proofs**

545 Question: For each theoretical result, does the paper provide the full set of assumptions and
546 a complete (and correct) proof?

547 Answer: [\[Yes\]](#)

548 Justification: The paper provides theoretical results with assumptions and proofs in Section
549 4

550 Guidelines:

- 551 • The answer [N/A] means that the paper does not include theoretical results.
- 552 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
553 referenced.
- 554 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 555 • The proofs can either appear in the main paper or the supplemental material, but if
556 they appear in the supplemental material, the authors are encouraged to provide a short
557 proof sketch to provide intuition.
- 558 • Inversely, any informal proof provided in the core of the paper should be complemented
559 by formal proofs provided in appendix or supplemental material.
- 560 • Theorems and Lemmas that the proof relies upon should be properly referenced.

561 4. Experimental result reproducibility

562 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
563 perimental results of the paper to the extent that it affects the main claims and/or conclusions
564 of the paper (regardless of whether the code and data are provided or not)?

565 Answer: [Yes],

566 Justification: The experiments can be reproduced by following experimental setup in Section
567 5

568 Guidelines:

- 569 • The answer [N/A] means that the paper does not include experiments.
- 570 • If the paper includes experiments, a [No] answer to this question will not be perceived
571 well by the reviewers: Making the paper reproducible is important, regardless of
572 whether the code and data are provided or not.
- 573 • If the contribution is a dataset and/or model, the authors should describe the steps taken
574 to make their results reproducible or verifiable.
- 575 • Depending on the contribution, reproducibility can be accomplished in various ways.
576 For example, if the contribution is a novel architecture, describing the architecture fully
577 might suffice, or if the contribution is a specific model and empirical evaluation, it may
578 be necessary to either make it possible for others to replicate the model with the same
579 dataset, or provide access to the model. In general, releasing code and data is often
580 one good way to accomplish this, but reproducibility can also be provided via detailed
581 instructions for how to replicate the results, access to a hosted model (e.g., in the case
582 of a large language model), releasing of a model checkpoint, or other means that are
583 appropriate to the research performed.
- 584 • While NeurIPS does not require releasing code, the conference does require all submis-
585 sions to provide some reasonable avenue for reproducibility, which may depend on the
586 nature of the contribution. For example
 - 587 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
588 to reproduce that algorithm.
 - 589 (b) If the contribution is primarily a new model architecture, the paper should describe
590 the architecture clearly and fully.
 - 591 (c) If the contribution is a new model (e.g., a large language model), then there should
592 either be a way to access this model for reproducing the results or a way to reproduce
593 the model (e.g., with an open-source dataset or instructions for how to construct
594 the dataset).
 - 595 (d) We recognize that reproducibility may be tricky in some cases, in which case
596 authors are welcome to describe the particular way they provide for reproducibility.
597 In the case of closed-source models, it may be that access to the model is limited in
598 some way (e.g., to registered users), but it should be possible for other researchers
599 to have some path to reproducing or verifying the results.

600 5. Open access to data and code

601 Question: Does the paper provide open access to the data and code, with sufficient instruc-
602 tions to faithfully reproduce the main experimental results, as described in supplemental
603 material?

604 Answer: [No]

605 Justification: The data is open-sourced but the code is not open-sourced

606 Guidelines:

- 607 • The answer [N/A] means that paper does not include experiments requiring code.
- 608 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
609 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 610 • While we encourage the release of code and data, we understand that this might not
611 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
612 including code, unless this is central to the contribution (e.g., for a new open-source
613 benchmark).
- 614 • The instructions should contain the exact command and environment needed to run to
615 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
616 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 617 • The authors should provide instructions on data access and preparation, including how
618 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 619 • The authors should provide scripts to reproduce all experimental results for the new
620 proposed method and baselines. If only a subset of experiments are reproducible, they
621 should state which ones are omitted from the script and why.
- 622 • At submission time, to preserve anonymity, the authors should release anonymized
623 versions (if applicable).
- 624 • Providing as much information as possible in supplemental material (appended to the
625 paper) is recommended, but including URLs to data and code is permitted.

626 6. Experimental setting/details

627 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
628 rameters, how they were chosen, type of optimizer) necessary to understand the results?

629 Answer: [Yes]

630 Justification: The experimental setting are clearly presented

631 Guidelines:

- 632 • The answer [N/A] means that the paper does not include experiments.
- 633 • The experimental setting should be presented in the core of the paper to a level of detail
634 that is necessary to appreciate the results and make sense of them.
- 635 • The full details can be provided either with the code, in appendix, or as supplemental
636 material.

637 7. Experiment statistical significance

638 Question: Does the paper report error bars suitably and correctly defined or other appropriate
639 information about the statistical significance of the experiments?

640 Answer: [No]

641 Justification: The paper follows most of related work to conduct experiments

642 Guidelines:

- 643 • The answer [N/A] means that the paper does not include experiments.
- 644 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
645 intervals, or statistical significance tests, at least for the experiments that support the
646 main claims of the paper.
- 647 • The factors of variability that the error bars are capturing should be clearly stated (for
648 example, train/test split, initialization, random drawing of some parameter, or overall
649 run with given experimental conditions).
- 650 • The method for calculating the error bars should be explained (closed form formula,
651 call to a library function, bootstrap, etc.)

- 652 • The assumptions made should be given (e.g., Normally distributed errors).
- 653 • It should be clear whether the error bar is the standard deviation or the standard error
- 654 of the mean.
- 655 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 656 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 657 of Normality of errors is not verified.
- 658 • For asymmetric distributions, the authors should be careful not to show in tables or
- 659 figures symmetric error bars that would yield results that are out of range (e.g., negative
- 660 error rates).
- 661 • If error bars are reported in tables or plots, the authors should explain in the text how
- 662 they were calculated and reference the corresponding figures or tables in the text.

663 8. Experiments compute resources

664 Question: For each experiment, does the paper provide sufficient information on the com-
 665 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 666 the experiments?

667 Answer: [Yes],

668 Justification: 8x H100 GPUs

669 Guidelines:

- 670 • The answer [N/A] means that the paper does not include experiments.
- 671 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 672 or cloud provider, including relevant memory and storage.
- 673 • The paper should provide the amount of compute required for each of the individual
- 674 experimental runs as well as estimate the total compute.
- 675 • The paper should disclose whether the full research project required more compute
- 676 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 677 didn't make it into the paper).

678 9. Code of ethics

679 Question: Does the research conducted in the paper conform, in every respect, with the
 680 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

681 Answer: [Yes]

682 Justification: The paper is conducted with the NeurIPS Code of Ethics

683 Guidelines:

- 684 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
- 685 Ethics.
- 686 • If the authors answer [No], they should explain the special circumstances that require a
- 687 deviation from the Code of Ethics.
- 688 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 689 eration due to laws or regulations in their jurisdiction).

690 10. Broader impacts

691 Question: Does the paper discuss both potential positive societal impacts and negative
 692 societal impacts of the work performed?

693 Answer: [N/A]

694 Justification: There is no societal impact of the work performed.

695 Guidelines:

- 696 • The answer [N/A] means that there is no societal impact of the work performed.
- 697 • If the authors answer [N/A] or [No], they should explain why their work has no societal
- 698 impact or why the paper does not address societal impact.
- 699 • Examples of negative societal impacts include potential malicious or unintended uses
- 700 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 701 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 702 groups), privacy considerations, and security considerations.

- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

718 11. Safeguards

719 Question: Does the paper describe safeguards that have been put in place for responsible
720 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
721 image generators, or scraped datasets)?

722 Answer: [N/A]

723 Justification: The paper poses no such risks.

724 Guidelines:

- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

735 12. Licenses for existing assets

736 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
737 the paper, properly credited and are the license and terms of use explicitly mentioned and
738 properly respected?

739 Answer: [Yes]

740 Justification: The assets are properly credited and the licenses are properly respected

741 Guidelines:

- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

755 • If this information is not available online, the authors are encouraged to reach out to
756 the asset’s creators.

757 **13. New assets**

758 Question: Are new assets introduced in the paper well documented and is the documentation
759 provided alongside the assets?

760 Answer: [N/A]

761 Justification: The paper does not release new assets.

762 Guidelines:

- 763 • The answer [N/A] means that the paper does not release new assets.
- 764 • Researchers should communicate the details of the dataset/code/model as part of their
765 submissions via structured templates. This includes details about training, license,
766 limitations, etc.
- 767 • The paper should discuss whether and how consent was obtained from people whose
768 asset is used.
- 769 • At submission time, remember to anonymize your assets (if applicable). You can either
770 create an anonymized URL or include an anonymized zip file.

771 **14. Crowdsourcing and research with human subjects**

772 Question: For crowdsourcing experiments and research with human subjects, does the paper
773 include the full text of instructions given to participants and screenshots, if applicable, as
774 well as details about compensation (if any)?

775 Answer: [N/A]

776 Justification: The paper does not involve crowdsourcing nor research with human subjects

777 Guidelines:

- 778 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
779 with human subjects.
- 780 • Including this information in the supplemental material is fine, but if the main contribu-
781 tion of the paper involves human subjects, then as much detail as possible should be
782 included in the main paper.
- 783 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
784 or other labor should be paid at least the minimum wage in the country of the data
785 collector.

786 **15. Institutional review board (IRB) approvals or equivalent for research with human
787 subjects**

788 Question: Does the paper describe potential risks incurred by study participants, whether
789 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
790 approvals (or an equivalent approval/review based on the requirements of your country or
791 institution) were obtained?

792 Answer: [N/A]

793 Justification: The paper does not involve crowdsourcing nor research with human subjects.

794 Guidelines:

- 795 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
796 with human subjects.
- 797 • Depending on the country in which research is conducted, IRB approval (or equivalent)
798 may be required for any human subjects research. If you obtained IRB approval, you
799 should clearly state this in the paper.
- 800 • We recognize that the procedures for this may vary significantly between institutions
801 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
802 guidelines for their institution.
- 803 • For initial submissions, do not include any information that would break anonymity (if
804 applicable), such as the institution conducting the review.

805 **16. Declaration of LLM usage**

806 Question: Does the paper describe the usage of LLMs if it is an important, original, or
807 non-standard component of the core methods in this research? Note that if the LLM is used
808 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
809 scientific rigor, or originality of the research, declaration is not required.

810 Answer: [N/A]

811 Justification: The core method development in this research does not involve LLMs as any
812 important, original, or non-standard components

- 813 • The answer [N/A] means that the core method development in this research does not
814 involve LLMs as any important, original, or non-standard components.
- 815 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
816 be described.