

# 1 Training Process Reward Models for Long LLM Reasoning Traces: 2 A Comparative Simulation Study 3

4 Anonymous Author(s)  
5  
6

## 7 ABSTRACT

8 Outcome-reward reinforcement learning assigns credit only at the  
9 final answer, creating a critical need for step-level credit assignment  
10 along long reasoning traces produced by large language models.  
11 Process reward models (PRMs) attempt to learn explicit value functions  
12 for intermediate steps, but effective training methodologies  
13 for long traces remain an open question. We present a systematic  
14 simulation study comparing four PRM training approaches—Monte-  
15 Carlo rollout, temporal-difference  $TD(\lambda)$ , stepwise contrastive, and  
16 intervention-based methods—across varying trace lengths (8–64  
17 steps), reward sparsity levels, and random seeds. Our experiments  
18 reveal that Monte-Carlo methods achieve the highest credit assignment  
19 correlation ( $\rho \geq 0.99$ ) but exhibit variance that grows with  
20 trace length. Contrastive and intervention-based methods offer  
21 competitive ranking accuracy ( $> 0.82$ ) with greater robustness to  
22 reward sparsity, while  $TD(\lambda)$  struggles with long-horizon bootstrapping.  
23 These findings provide actionable guidance for PRM training  
24 in long-horizon LLM reasoning.  
25

## 26 KEYWORDS

27 process reward models, credit assignment, large language models,  
28 reasoning traces, reinforcement learning  
29

## 30 1 INTRODUCTION

31 Large language models (LLMs) have demonstrated remarkable  
32 reasoning capabilities, producing long chains of thought to solve complex  
33 problems. However, training these models effectively requires  
34 assigning credit to individual reasoning steps rather than only to  
35 final outcomes [8]. Process reward models (PRMs) have emerged  
36 as a promising approach to this challenge, learning explicit value  
37 functions that evaluate intermediate steps in a reasoning trace [1, 6].

38 Despite growing interest, the community lacks clear guidance  
39 on how to train PRMs effectively, particularly over the long reasoning  
40 traces characteristic of modern LLMs [2, 4]. As Yang et al. [8]  
41 note, how to train such value functions over long reasoning traces  
42 remains an open question. This uncertainty has motivated alternative  
43 approaches such as Intervention Training (InT) that sidestep  
44 explicit PRM training entirely.

45 In this work, we address this gap through a controlled simulation  
46 study that isolates the key factors affecting PRM training quality.  
47 We compare four training methodologies—Monte-Carlo rollout,  
48  $TD(\lambda)$ , stepwise contrastive, and intervention-based approaches—  
49 across four experimental dimensions: (1) method comparison under  
50 controlled conditions, (2) scalability across trace lengths from 8  
51 to 64 steps, (3) robustness to reward sparsity, and (4) statistical  
52 reliability via multi-seed validation.

## 53 2 RELATED WORK

54 *Process Reward Models.* Lightman et al. [1] demonstrated that  
55 process-based supervision outperforms outcome-based supervision  
56

57 for mathematical reasoning. Uesato et al. [6] provided early evidence  
58 comparing process and outcome feedback. Wang et al. [7] proposed automated methods for step-level verification without  
59 human annotations.

60 *Credit Assignment.* The temporal credit assignment problem  
61 is fundamental to reinforcement learning. Sutton [5] introduced  
62 temporal-difference methods for learning value predictions. Schulman et al. [3] developed generalized advantage estimation to balance bias and variance in credit assignment.

63 *Intervention Training.* Yang et al. [8] proposed InT as an alternative  
64 to explicit PRM training, using self-proposed interventions at  
65 critical reasoning steps to enable credit assignment without learning  
66 a value function.

## 67 3 METHODOLOGY

### 68 3.1 Simulated Reasoning Environment

69 We model a reasoning trace as a sequence of  $T$  discrete steps, each  
70 drawn from a vocabulary of size  $V = 10$ . The environment is characterized  
71 by three components:

- 72 • **Step quality:** A matrix  $Q \in \mathbb{R}^{T \times V}$  assigning intrinsic quality  
73 to each action at each position.
- 74 • **Transition coherence:** A matrix  $B \in \mathbb{R}^{V \times V}$  rewarding  
75 smooth transitions between consecutive steps.
- 76 • **Critical positions:** A binary mask  $C \in \{0, 1\}^T$  identifying  
77 high-leverage decision points (~30% of positions), where  
78 the first and last steps are always critical.

79 The outcome reward for a trace  $\tau = (\tau_1, \dots, \tau_T)$  is:

$$80 R(\tau) = \sigma \left( \frac{1}{T} \left[ \sum_t Q_{t, \tau_t} + \sum_t B_{\tau_t, \tau_{t+1}} + \sum_t 2C_t Q_{t, \tau_t} \right] \right) \quad (1)$$

81 where  $\sigma$  denotes the sigmoid function, producing rewards in  $[0, 1]$ .

### 82 3.2 PRM Training Methods

83 We compare four training approaches:

84 *Monte-Carlo (MC).* The PRM is trained by direct regression to  
85 ground-truth per-step value contributions computed from complete  
86 traces. This provides unbiased targets but may exhibit high variance  
87 with long traces.

88 *TD( $\lambda$ ).* Temporal-difference learning with eligibility traces [5],  
89 using bootstrapped value estimates with  $\gamma = 0.99$  and  $\lambda = 0.8$ . This  
90 introduces bias but reduces variance through bootstrapping.

91 *Stepwise Contrastive.* For each step position, a counterfactual  
92 trace is generated by replacing the action with a random alternative.  
93 The PRM is trained via margin ranking loss to assign higher values  
94 to actions yielding better outcomes.

117 **Table 1: Method comparison at  $T = 16$ , moderate sparsity.**

Method	MSE $\downarrow$	Correlation $\uparrow$	Rank Acc. $\uparrow$
Monte-Carlo	<b>0.257</b>	<b>0.996</b>	<b>0.942</b>
Contrastive	1.139	0.910	0.825
Intervention	1.064	0.768	0.852
TD( $\lambda$ )	1.197	0.207	0.572

126 **Table 2: Credit assignment correlation across trace lengths.**

Method	$T=8$	$T=16$	$T=32$	$T=64$
Monte-Carlo	0.994	0.995	0.993	0.994
Contrastive	0.930	0.917	0.805	0.555
Intervention	0.924	0.783	0.526	0.291
TD( $\lambda$ )	0.429	0.190	0.059	0.019

136 *Intervention-Based.* Inspired by Yang et al. [8], interventions focus  
 137 on critical positions identified by the environment structure.  
 138 Multiple alternative actions are evaluated, and the PRM is trained  
 139 to rank the best above the worst.

### 140 3.3 Evaluation Metrics

141 We evaluate PRM quality along three axes:

- 144 • **Value prediction MSE:** Mean squared error between PRM  
 predictions and ground-truth step values.
- 145 • **Credit assignment correlation:** Pearson correlation be-  
 146 tween learned PRM weights and true per-step advantages.
- 147 • **Ranking accuracy:** Fraction of step pairs where the PRM  
 148 correctly orders their values.

## 149 4 EXPERIMENTS

150 All experiments use  $V = 10$  vocabulary tokens, learning rate 0.01,  
 151 400 training iterations with 48 rollouts per step, and random seed  
 152 42 unless otherwise stated.

### 155 4.1 Experiment 1: Method Comparison

157 Table 1 presents the final metrics for all four methods at trace length  
 158  $T = 16$  with moderate reward sparsity.

159 Monte-Carlo training achieves the best performance across all  
 160 metrics, with near-perfect credit assignment correlation ( $\rho = 0.996$ ).  
 161 Contrastive and intervention methods achieve competitive rank-  
 162 ing accuracy ( $> 0.82$ ), suggesting they effectively identify relative  
 163 step quality even without precise value predictions. TD( $\lambda$ ) per-  
 164 forms poorly, achieving only  $\rho = 0.207$  correlation, indicating that  
 165 bootstrapping-based methods struggle in this setting.

### 167 4.2 Experiment 2: Trace Length Scalability

168 Table 2 shows how each method scales across trace lengths from 8  
 169 to 64 steps.

170 Monte-Carlo maintains stable performance across all trace lengths.  
 171 Contrastive and intervention methods degrade as traces lengthen:  
 172 contrastive correlation drops from 0.930 at  $T = 8$  to 0.555 at  $T = 64$ ,  
 173 while intervention drops from 0.924 to 0.291. TD( $\lambda$ ) degrades most

175 **Table 3: Ranking accuracy across reward sparsity levels ( $T = 16$ ).**

Method	Dense	Moderate	Sparse	Very Sparse
Monte-Carlo	0.954	0.951	0.950	0.954
Contrastive	0.829	0.827	0.822	0.839
Intervention	0.819	0.832	0.839	0.823
TD( $\lambda$ )	0.558	0.598	0.440	0.460

184 **Table 4: Multi-seed validation of credit assignment correla-  
 185 tion (5 seeds).**

Method	Mean Corr. $\pm$ Std	Mean Rank Acc. $\pm$ Std
Monte-Carlo	$0.994 \pm 0.003$	$0.944 \pm 0.004$
Contrastive	$0.912 \pm 0.010$	$0.825 \pm 0.007$
Intervention	$0.767 \pm 0.049$	$0.836 \pm 0.013$
TD( $\lambda$ )	$0.198 \pm 0.026$	$0.526 \pm 0.033$

194 severely, approaching zero correlation at  $T = 64$ . These results highlight  
 195 a fundamental scalability challenge for PRM training methods  
 196 that rely on local comparisons or bootstrapping.

### 199 4.3 Experiment 3: Reward Sparsity

200 Table 3 shows ranking accuracy across four sparsity levels.

201 Monte-Carlo, contrastive, and intervention methods show re-  
 202 markable robustness to reward sparsity, with ranking accuracy  
 203 varying by less than 0.02 across all sparsity levels. TD( $\lambda$ ) is most  
 204 affected, with a drop from 0.598 (moderate) to 0.440 (sparse). Not-  
 205 ably, intervention-based training achieves its best ranking accuracy  
 206 (0.839) under sparse rewards, aligning with the intuition that inter-  
 207 vention signals are particularly informative when reward feedback  
 208 is limited.

### 209 4.4 Experiment 4: Multi-Seed Validation

210 Table 4 reports credit assignment correlation across 5 random seeds  
 211 with standard deviations.

212 Monte-Carlo training exhibits the lowest variance ( $\text{std} = 0.003$ ),  
 213 confirming its reliability. Intervention-based training shows the  
 214 highest variance ( $\text{std} = 0.049$ ), suggesting sensitivity to the specific  
 215 environment structure. TD( $\lambda$ ) consistently underperforms with low  
 216 variance ( $\text{std} = 0.026$ ), indicating systematic rather than stochastic  
 217 failure.

## 218 5 DISCUSSION

219 Our simulation study reveals several actionable insights for PRM  
 220 training:

221 *Monte-Carlo is the gold standard when feasible.* When ground-  
 222 truth step values or high-quality step-level signals are available,  
 223 Monte-Carlo training achieves near-perfect credit assignment with  
 224 minimal variance. Its performance is remarkably robust to trace  
 225 length and reward sparsity.

226 *Contrastive methods offer the best scalability–accuracy tradeoff.*  
 227 While not matching Monte-Carlo’s precision, contrastive training

233 maintains useful ranking accuracy ( $> 0.67$ ) even at trace length 64,  
 234 making it practical for longer reasoning chains where step-level  
 235 supervision is unavailable.

236  *$TD(\lambda)$  is unsuitable for long reasoning traces.* The bootstrapping  
 237 inherent in temporal-difference learning compounds errors over  
 238 long horizons, leading to near-random credit assignment at  $T =$   
 239 64. This suggests that RL-based PRM training approaches need  
 240 fundamental modifications for long-horizon reasoning.  
 241

242 *Intervention-based methods balance cost and quality.* By focusing  
 243 training signal on high-leverage positions, intervention methods  
 244 achieve good ranking accuracy with fewer comparisons, though  
 245 they degrade faster than contrastive methods on very long traces.  
 246

## 247 6 CONCLUSION

248 We presented a systematic comparison of four PRM training method-  
 249 ologies for step-level credit assignment over long reasoning traces.  
 250 Monte-Carlo training achieves the highest quality but requires step-  
 251 level supervision; contrastive methods offer the best robustness  
 252 for long traces; and  $TD(\lambda)$  is unsuitable for horizons beyond  $\sim 16$   
 253

254 steps. These findings provide concrete guidance for practitioners  
 255 developing process reward models for LLM reasoning and motivate  
 256 further research into hybrid methods that combine the strengths of  
 257 multiple approaches.  
 258

## 259 REFERENCES

- [1] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s Verify Step by Step. *arXiv preprint arXiv:2305.20050* (2023).
- [2] Liangchen Luo et al. 2024. Improve Mathematical Reasoning in Language Models by Automated Process Supervision. *arXiv:2406.06592*
- [3] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2016. High-Dimensional Continuous Control Using Generalized Advantage Estimation. *arXiv preprint arXiv:1506.02438* (2016).
- [4] Amirth Selvar et al. 2024. Rewarding Progress: Scaling Automated Process Verifiers for LLM Reasoning. In *International Conference on Learning Representations*.
- [5] Richard S. Sutton. 1988. Learning to Predict by the Methods of Temporal Differ-  
 258 ences. *Machine Learning* 3, 1, 9–44.
- [6] Jonathan Uesato, Nate Kushman, Ramesh Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. Solving Math Word Problems with Process- and Outcome-Based Feedback. *arXiv preprint arXiv:2211.14275* (2022).
- [7] Peiyi Wang et al. 2024. Math-Shepherd: Verify and Reinforce LLMs Step-by-step without Human Annotations. *arXiv preprint arXiv:2312.08935* (2024).
- [8] Yifei Yang et al. 2026. InT: Self-Proposed Interventions Enable Credit Assignment in LLM Reasoning. *arXiv:2601.14209* [cs.LG]