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

Research

## ABSTRACT

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

## KEYWORDS

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

## 1 INTRODUCTION

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

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

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

## 2 RELATED WORK

*Process Reward Models.* Lightman et al. [1] demonstrated that process-based supervision outperforms outcome-based supervision for mathematical reasoning. Uesato et al. [6] provided early evidence comparing process and outcome feedback. Wang et al. [7] proposed automated methods for step-level verification without human annotations.

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

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

## 3 METHODOLOGY

### 3.1 Simulated Reasoning Environment

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

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

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

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

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

### 3.2 PRM Training Methods

We compare four training approaches:

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

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

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

| Method | MSE ↓ | Correlation ↑ | Rank Acc. ↑ |
|---|---|---|---|
| Monte-Carlo | **0.257** | **0.996** | **0.942** |
| Contrastive | 1.139 | 0.910 | 0.825 |
| Intervention | 1.064 | 0.768 | 0.852 |
| TD($\lambda$) | 1.197 | 0.207 | 0.572 |

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

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

### 3.3 Evaluation Metrics

We evaluate PRM quality along three axes:

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

## 4 EXPERIMENTS

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

### 4.1 Experiment 1: Method Comparison

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

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

### 4.2 Experiment 2: Trace Length Scalability

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

Monte-Carlo maintains stable performance across all trace lengths. Contrastive and intervention methods degrade as traces lengthen: contrastive correlation drops from 0.930 at $T = 8$ to 0.555 at $T = 64$, while intervention drops from 0.924 to 0.291. TD($\lambda$) degrades most severely, approaching zero correlation at $T = 64$. These results highlight a fundamental scalability challenge for PRM training methods that rely on local comparisons or bootstrapping.

**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 |

**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 |

**Table 4: Multi-seed validation of credit assignment correlation (5 seeds).**

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

### 4.3 Experiment 3: Reward Sparsity

Table 3 shows ranking accuracy across four sparsity levels.

Monte-Carlo, contrastive, and intervention methods show remarkable robustness to reward sparsity, with ranking accuracy varying by less than 0.02 across all sparsity levels. TD($\lambda$) is most affected, with a drop from 0.598 (moderate) to 0.440 (sparse). Notably, intervention-based training achieves its best ranking accuracy (0.839) under sparse rewards, aligning with the intuition that intervention signals are particularly informative when reward feedback is limited.

### 4.4 Experiment 4: Multi-Seed Validation

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

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

## 5 DISCUSSION

Our simulation study reveals several actionable insights for PRM training:

*Monte-Carlo is the gold standard when feasible.* When ground-truth step values or high-quality step-level signals are available,

Monte-Carlo training achieves near-perfect credit assignment with minimal variance. Its performance is remarkably robust to trace length and reward sparsity.

*Contrastive methods offer the best scalability–accuracy tradeoff.* While not matching Monte-Carlo's precision, contrastive training maintains useful ranking accuracy ($> 0.67$) even at trace length 64, making it practical for longer reasoning chains where step-level supervision is unavailable.

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

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

## 6 CONCLUSION

We presented a systematic comparison of four PRM training methodologies for step-level credit assignment over long reasoning traces.

Monte-Carlo training achieves the highest quality but requires step-level supervision; contrastive methods offer the best robustness for long traces; and TD($\lambda$) is unsuitable for horizons beyond ~16 steps. These findings provide concrete guidance for practitioners developing process reward models for LLM reasoning and motivate further research into hybrid methods that combine the strengths of multiple approaches.

## 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] Amrith Setlur 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 Differences. *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]