

# PROCESS SUPERVISION-GUIDED POLICY OPTIMIZATION FOR CODE GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement learning (RL) with unit test feedback has enhanced large language models' (LLMs) code generation, but relies on sparse rewards provided only after complete code evaluation, limiting learning efficiency and incremental improvements. When generated code fails all unit tests, no learning signal is received, hindering progress on complex tasks. To address this, we propose a Process Reward Model (PRM) that delivers dense, line-level feedback on code correctness during generation, mimicking human code refinement and providing immediate guidance. We explore various strategies for training PRMs and integrating them into the RL framework, finding that using PRMs both as dense rewards and for value function initialization significantly boosts performance. Our experimental results also highlight the effectiveness of PRMs in enhancing RL-driven code generation, especially for long-horizon scenarios.

## 1 INTRODUCTION

The rapid advancement of large language models (LLMs) has revolutionized code generation, enabling models to achieve near-human performance on programming tasks (Chen et al., 2021a; Li et al., 2022; OpenAI, 2023). These models have demonstrated remarkable abilities to generate syntactically correct and functionally viable code snippets, significantly aiding software development processes. Building upon these successes, recent research has explored the use of reinforcement learning (RL) from unit test feedback to further enhance the code generation capabilities of LLMs (Le et al., 2022; Shojaaee et al., 2023; Liu et al., 2023; Dou et al., 2024). By incorporating unit tests as a reward signal, these methods aim to guide LLMs toward generating code that not only compiles but also passes specified test cases, thereby improving overall code reliability and quality.

However, a significant challenge arises from the nature of the reward signals derived from unit tests. These signals are inherently sparse, as they are only received at the end of an episode after the entire code snippet has been generated and evaluated. This delay in feedback impedes learning efficiency and limits the model's ability to make incremental improvements during code generation. When an LLM fails to generate code that passes any unit tests, it receives no meaningful learning signal, making it difficult to learn to solve more complex coding problems. In contrast, human programmers typically do not rewrite code from scratch when their programs fail unit tests. Instead, they analyze the code to pinpoint and fix errors, leveraging their understanding of programming logic and structure to iteratively improve upon the current version. This process of step-by-step refinement, which involves receiving and acting upon fine-grained feedback, is missing in the current RL training loop for code generation from unit test feedback.

To address this limitation, we propose integrating a Process Reward Model (PRM) (Lightman et al., 2023; Wang et al., 2024a) into the RL training framework for code generation. A PRM provides dense signals by offering line-level feedback that indicates the correctness of each generated line of code. This fine-grained feedback mechanism mimics the human approach to code refinement and has the potential to enhance learning efficiency by providing immediate guidance during code generation. While the concept of using PRMs is intuitive, training an effective PRM and integrating it into RL training is non-trivial. Challenges include accurately modeling the correctness of partial code snippets and ensuring stable and effective training dynamics when combining PRM-generated signals with traditional RL methods. Although previous research has attempted to incorporate PRMs into LLM RL training (Wang et al., 2024a), these efforts have been limited to the mathematical domain and have not fully explored the complexities involved.

In this work, we conduct a comprehensive analysis of how PRMs can be integrated into RL training for code generation. We explore various strategies for training a robust code PRM and investigate different methods of utilizing it to improve code generation performance. Based on our experiments, we provide a practical recipe for successfully using PRMs and integrating them into RL training in the context of code generation. Notably, one of our key findings is that using PRMs concurrently as both dense rewards and value function initialization in RL training improves learning efficiency and leads to a significant performance improvement. Our contributions can be summarized as follows:

- We present a practical and effective pipeline that automatically generates process-level supervision data, trains a PRM based on these supervisions, and integrates it into RL training to provide dense feedback signals. To the best of our knowledge, we are the first to demonstrate that PRMs can benefit RL from unit test feedback in code generation.
- We systematically investigate how to best integrate PRMs into RL training, exploring different strategies for training code PRMs and utilizing them to enhance code generation. Our findings are distilled into a practical guideline for effectively incorporating PRMs in RL for code generation.
- Following our proposed methodology, we demonstrate significant improvements in pass rates across HumanEval, MBPP, and LiveCodeBench benchmarks for two baseline LLMs. Additionally, we highlight key insights: (1) the synergy between dense rewards and value initialization maximizes PRM effectiveness, (2) PRMs encourage exploration and improve learning efficiency, and (3) PRMs enhance code generation, particularly in long-horizon scenarios.

## 2 PROBLEM FORMALIZATION

In code generation tasks, we define a code generation problem as a sequence of tokens  $\mathbf{x} = (x_1, x_2, \dots, x_m)$ , where each  $x_i$  denotes the  $i$ -th element or token of the input prompt, which may include problem descriptions. The primary objective for the model in this context is to process the given input  $\mathbf{x}$  and generate a coherent and syntactically correct sequence of code tokens. This sequence is denoted as  $\mathbf{y} = (\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(T)})$ , where  $T$  represents the total number of code generation steps. Each individual code generation step,  $\mathbf{y}^{(t)}$ ,  $t = 1, 2, \dots, T$ , is composed of a series of tokens  $y_1^{(t)}, y_2^{(t)}, \dots, y_{n_t}^{(t)}$ , where  $y_i^{(t)}$  corresponds to the  $i$ -th token within the  $t$ -th step, and  $n_t$  denotes the number of tokens in this step. Typically, a pre-trained language model (LM), denoted as  $p_\theta$ , is employed to model the conditional probability distribution of the code generation steps  $\mathbf{y}$ , given the code generation problem  $\mathbf{x}$ , which is mathematically represented as  $p_\theta(\mathbf{y} | \mathbf{x})$ , parameterized by  $\theta$ . The model is optimized through training on a dataset  $\mathcal{D}_{\mathbf{x}\mathbf{y}}$  containing pairs of prompts and their corresponding code solutions. This training process, often referred to as Supervised Fine-Tuning (SFT), involves maximizing the log-likelihood of the dataset.

### 2.1 BASELINE METHOD: REINFORCEMENT LEARNING FROM UNIT TEST FEEDBACK

Code generation tasks can be formulated within a Reinforcement Learning (RL) framework, where code generation is treated as a sequence of decision-making steps. Once the model has undergone SFT, the RL phase is employed to refine the model’s ability to generate functionally correct code using feedback from unit tests (Liu et al., 2023). Unit test feedback is derived by executing the generated program on predefined test cases. The feedback serves as a signal for learning and can be transformed into a reward. A simple reward function based on the outcome of the unit tests could be defined as follows:

$$R_{\text{UT}}(\mathbf{x}, \mathbf{y}) = \begin{cases} 1, & \text{if the program } \mathbf{y} \text{ passes all unit tests} \\ 0, & \text{otherwise} \end{cases}$$

This binary reward formulation encourages the model to generate programs that can successfully pass all unit test cases. Given a collection of unlabeled code generation prompts  $\mathcal{D}_{\mathbf{x}}$ , the model  $p_\theta$  is optimized to maximize the expected reward over all possible code generation trajectories.

## 3 PROCESS SUPERVISION-GUIDED POLICY OPTIMIZATION

While the Reinforcement Learning from Unit Test Feedback (RLTF) offers a framework for improving code generation models, it suffers from significant limitations due to the sparsity of its reward signal. The binary nature of unit test feedback—indicating only whether the entire program passes

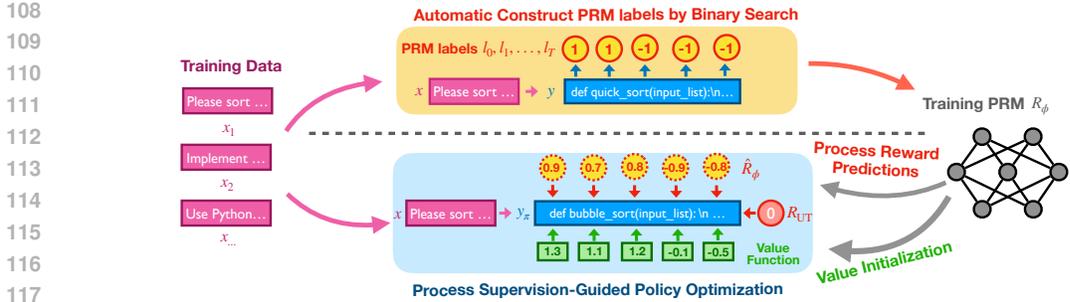


Figure 1: Overview of our method. Our approach consists of two main components: (1) a binary search-based method for automating PRM training data labeling, used to train a code PRM; and (2) the integration of PRM into RL training, where it serves as (a) the initialization for the value model and (b) an evaluator assessing the correctness of each line of code, providing dense reward signals.

or fails—provides no guidance on which specific parts of the code contributed to the outcome. This lack of intermediate feedback makes it challenging for the model to identify and correct errors during training, leading to slow convergence and suboptimal performance. In contrast, human programmers iteratively develop and refine their code. When a program fails to pass unit tests, they do not typically rewrite it from scratch. Instead, they analyze the code to pinpoint and fix errors, leveraging their understanding of programming logic and structure. This process of step-by-step refinement is crucial for efficient problem-solving.

Motivated by this observation, we propose **Process Supervision-Guided Policy Optimization (PS-GPO)**, a method that integrates fine-grained feedback into the RL framework. Figure 1 illustrates the overview of our approach. By providing intermediate rewards that assess the correctness of partial code sequences, our approach guides the model more effectively toward generating functionally correct programs. This is achieved using a Process Reward Model (PRM) (Lightman et al., 2023), which serves as 1) the initialization for the value model and 2) an evaluator for each code generation step, providing dense reward signals that address the limitations of end-of-trajectory rewards.

### 3.1 PROCESS SUPERVISION VIA PROCESS REWARD MODELS

Our method introduces a PRM to assess the correctness of each **line** of the code during the generation process. The PRM serves as an oracle that provides intermediate rewards based on the potential of the current code prefix to be extended into a correct program. By offering intermediate feedback, the PRM helps the model identify and reinforce beneficial code generation patterns while discouraging those that introduce errors. This fine-grained feedback mirrors the human approach to coding, where programmers continuously evaluate and adjust their code.

#### 3.1.1 DATA COLLECTION

To effectively train the PRM, we require a dataset that provides fine-grained annotations indicating the correctness of partial code sequences. However, manually annotating the correctness of each line of code generated by the model is costly and not scalable. Instead, we employ an automated approach inspired by techniques used in recent works (Wang et al., 2024a;b; Luo et al., 2024). Our method leverages the model’s own capabilities to generate completions for partial code prefixes and uses automated testing to assess their correctness. The key idea is to determine whether a partial code prefix can be extended—by any means—into a complete program that passes all unit tests. If so, we consider the prefix as potentially correct; otherwise, it is labeled as incorrect.

Given a prompt  $\mathbf{x}$ , we generate a complete code response  $\mathbf{y} = (\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(T)})$  using the current policy  $p_\theta$ . Our goal is to determine the correctness of each partial code prefix  $\mathbf{y}^{\leq t}$  for  $t = 1, 2, \dots, T$ . To achieve this, we employ a *best-of- $K$*  sampling strategy to approximate an oracle capable of completing partial code prefixes. For each partial code prefix  $\mathbf{y}^{\leq t}$ , we generate  $K$  potential completions  $\{c_k\}_{k=1}^K$  using the current policy. We then form full programs  $\mathcal{P}_k = (\mathbf{y}^{\leq t}, c_k)$  and execute them against the unit tests  $\mathcal{U}$ . If any of these programs pass all unit tests, we label the partial code prefix as **correct**; otherwise, it is labeled as **incorrect**.

To efficiently identify the transition point where errors occur, we employ a binary search over the code generation steps (Luo et al., 2024), which is formalized in Algorithm 1. For example, consider a code response divided into five steps ( $T = 5$ ), as shown in Figure 2. The partial prefix up to  $\mathbf{y}^{(3)}$  can be completed to pass all unit tests, so it is labeled as correct. The prefix up to  $\mathbf{y}^{(4)}$  cannot, meaning steps beyond  $\mathbf{y}^{(3)}$  are labeled as incorrect.

For each partial code prefix  $\mathbf{y}^{\leq m}$ , the label  $l_m$  is assigned based on the outcome of the completion attempts:

$$l_m = \begin{cases} +1, & \text{if any } \mathcal{P}_k \text{ passes all unit tests} \\ -1, & \text{otherwise} \end{cases} \quad (1)$$

which indicate whether the prefix is potentially correct (can be completed to a correct program) or incorrect (contains unrecoverable errors).

### 3.1.2 TRAINING

Using the collected data  $\{(\mathbf{x}, \mathbf{y}^{\leq m}, l_m)\}$ , we train the PRM  $R_\phi$  to predict the correctness of partial code prefixes. The PRM learns to assign higher rewards to prefixes labeled as correct and lower rewards to those labeled as incorrect. The training objective, i.e., Mean Squared Error (MSE), minimizes the discrepancy between the PRM’s predictions and the assigned labels:

$$\min_{\phi} \sum_{(\mathbf{x}, \mathbf{y}^{\leq m})} \left( R_\phi(\mathbf{x}, \mathbf{y}^{\leq m}) - l_m \right)^2 \quad (2)$$

This regression formulation allows the PRM to estimate the likelihood that a given prefix can be successfully completed. Notably, we also experimented with Cross-Entropy loss and empirically found that MSE loss yielded better performance.

## 3.2 INTEGRATING PRM INTO RL TRAINING

Given a learned PRM, we aim to identify best practices for enhancing code generation during RL training. While prior work has used PRMs to verify intermediate steps in mathematical tasks (Lightman et al., 2023; Wang et al., 2024a; Jiao et al., 2024; Wang et al., 2024b; Luo et al., 2024), their potential for guiding code generation remains largely unexplored. In mathematical domains, LLMs may generate correct answers with faulty reasoning (Lightman et al., 2023), making intermediate verification essential. However, in code generation, problems are typically accompanied by multiple unit tests, making it improbable for incorrect code to pass all tests. As a result, the emphasis on intermediate verification is less applicable. Instead, we propose leveraging PRMs as auxiliary sources of dense signals to facilitate better exploration during RL training. While preliminary attempts have been made to incorporate PRMs into RL training (Wang et al., 2024a), these efforts are limited and warrant a more thorough investigation. We explore the following methods to integrate PRMs effectively:

**PRM as Dense Rewards.** Similar to Wang et al. (2024a), we use PRMs to provide step-level reward signals that guide more efficient policy exploration during RL training. By rating the correctness of each line in the code response, the PRM supplies “dense” rewards that encourage the policy to explore more promising code paths.

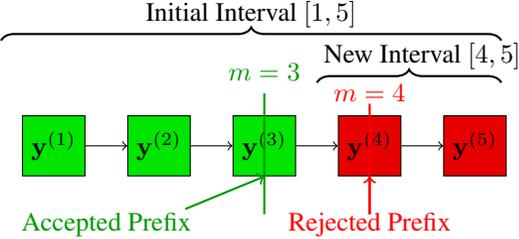


Figure 2: Binary search over code steps at line level to label prefixes. The first midpoint at  $m = 3$  is accepted, so the search interval moves to  $[4, 5]$ . The next midpoint at  $m = 4$  is rejected, indicating unrecoverable errors occur after step 3.

---

### Algorithm 1: Binary Search for Labeling Prefixes

---

**Require:** Prompt  $\mathbf{x}$ , response  $\mathbf{y}$ , policy  $p_\theta$ , unit tests  $\mathcal{U}$ , number of completions  $K$   
**Ensure:** Labels  $l_t$  for each prefix  $\mathbf{y}^{\leq t}$

- 1: Initialize lower bound  $L \leftarrow 1$ , upper bound  $R \leftarrow T$ , failure point  $F \leftarrow T + 1$
- 2: **while**  $L \leq R$  **do**
- 3:   Compute midpoint  $m \leftarrow \lfloor (L + R)/2 \rfloor$
- 4:   Set success flag  $S \leftarrow \text{False}$
- 5:   **for**  $k = 1$  to  $K$  **do**
- 6:     Generate completion  $\mathbf{c}_k \sim p_\theta(\cdot \mid \mathbf{y}^{\leq m})$
- 7:     Form full program  $\mathcal{P}_k \leftarrow (\mathbf{y}^{\leq m}, \mathbf{c}_k)$
- 8:     Execute  $\mathcal{P}_k$  with unit tests  $\mathcal{U}$
- 9:     **if**  $\mathcal{P}_k$  passes all unit tests **then**
- 10:       Set  $S \leftarrow \text{True}$ ; **break**
- 11:     **end if**
- 12:   **end for**
- 13:    $L \leftarrow$  if  $S = \text{True}$  then  $m + 1$  else  $L$
- 14:    $F \leftarrow$  if  $S = \text{True}$  then  $F$  else  $m$
- 15:    $R \leftarrow$  if  $S = \text{True}$  then  $R$  else  $m - 1$
- 16: **end while**
- 17: **for**  $t = 1$  to  $T$  **do**
- 18:    $l_t \leftarrow$  if  $t < F$  then  $+1$  else  $-1$
- 19: **end for**

---

**PRM as Value Initialization.** The PRM’s method of annotating code, by fixing a prefix  $\mathbf{y}^{\leq t}$  and rolling out the policy to sample correct responses, can be viewed as a “hard” value estimation of  $\mathbf{y}^{\leq t}$ . We hypothesize that the PRM’s capability to provide line-level feedback could serve as a useful inductive bias for initializing the value function in RL algorithms, which can ease the credit assignment problem by offering a more informed starting point.

## 4 EXPERIMENTAL RESULTS

### 4.1 EXPERIMENTAL SETUP

**Datasets and Evaluation.** We utilize in-house datasets to train our model for code generation. Specifically, the training set,  $\mathcal{D}_{\text{train}}$ , is a comprehensive Reinforcement Learning with Human Feedback (RLHF) dataset that includes, as a subset, approximately 30,000 diverse coding problems. Each of these problems is paired with unit tests designed to validate the functional correctness of the generated code. For evaluation, we use three widely adopted benchmarks: **HumanEval** (Chen et al., 2021b), which contains 164 hand-crafted programming problems; **MBPP** (Austin et al., 2021), a dataset of 974 crowd-sourced Python programming problems (where problems with IDs 11–510 are used for evaluation); and **LiveCodeBench** (Jain et al., 2024), a comprehensive benchmark designed to assess the code generation capabilities of LLMs. Among the various releases of LiveCodeBench, we use **LiveCodeBench v3**, which includes 612 coding tasks collected between May 2023 and July 2024. In our experiments, we restricted LLMs to a single-turn chat completion setting. For each coding problem, we directly input the problem to the model without using few-shot prompting. We generate 10 candidate responses for each problem, using a temperature of 0.2, nucleus sampling with top- $p=0.95$ , and top- $k$  sampling with  $k=128$ , following common practice. We adopt Pass@1 as the evaluation metric, in line with previous work (Kulal et al., 2019; Chen et al., 2021a; Jain et al., 2024).

**Base Models.** In our experiments, we employ two different base models, Qwen2.5-7B and InHouse-Lite, to evaluate the effectiveness of our proposed method. Qwen2.5-7B (Qwen-Team, 2024) is a publicly released causal language model from the Qwen series (Yang et al., 2024), recognized for its strong performance across diverse domains. InHouse-Lite is an in-house model of comparable size to Qwen2.5-7B but utilizes a different network architecture.

**SFT and RL Baseline.** Initially, both Qwen2.5-7B and InHouse-Lite are fine-tuned on our Supervised Fine-Tuning (SFT) dataset, yielding Qwen2.5-7B-SFT and InHouse-Lite-SFT, which serve as the starting points for the subsequent RLHF training phase. We then further optimize these SFT models ( $\pi_{\text{ref}}$ ) on the RLHF dataset  $\mathcal{D}_{\text{train}}$  using Proximal Policy Optimization (PPO) (Schulman et al., 2017), resulting in the RL models Qwen2.5-7B-RL and InHouse-Lite-RL ( $\pi_{\theta}$ ). In our setup, two types of Outcome Reward Models (ORMs) are employed as the objective functions for RL training. For non-coding prompts, we use a general reward model,  $R_{\text{general}}(\mathbf{x}, \mathbf{y})$ , derived from preference learning on a human-annotated dataset (Ouyang et al., 2022). For coding prompts, the ORM is defined as a binary indicator of whether the response passes all unit tests,  $R_{\text{UT}}(\mathbf{x}, \mathbf{y})$ . Following Ouyang et al. (2022), RLHF optimization objective is defined as:

$$\max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}_{\text{train}}} \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x})} [R(\mathbf{x}, \mathbf{y}) - \beta \text{KL}(\pi_{\theta} \parallel \pi_{\text{ref}})],$$

with  $R(\mathbf{x}, \mathbf{y}) = R_{\text{general}}(\mathbf{x}, \mathbf{y})$  for non-coding prompts and  $R(\mathbf{x}, \mathbf{y}) = R_{\text{UT}}(\mathbf{x}, \mathbf{y})$  for coding prompts.

**PRM Training.** To ensure that the PRM training data effectively covers the state space the language model may encounter during the next RL training phase, we sample policy models from various stages of the RL baseline training. Specifically, we select 4 checkpoints evenly spaced throughout the RL baseline model’s training process. For each checkpoint, we sample  $n$  responses for each coding prompt in the training dataset  $\mathcal{D}_{\text{train}}$ . For each sampled response, we apply the binary search labeling procedure described in Algorithm 1, using  $K = 20$  completions for each partial code prefix. The data collected from all checkpoints is then aggregated into a PRM training set, denoted as  $\mathcal{D}_{\text{PRM}}$ . We initialize the PRM with the value model from the RL baseline and fine-tune it on the aggregated dataset,  $\mathcal{D}_{\text{PRM}}$ , using the objective function defined in Eq. (2).

**Integrating PRM into RL.** As described in Section 3.2, we explore two methods for integrating the Process Reward Model (PRM) into RL training: 1) using PRM as a source of dense reward signals (**DenseReward**) and 2) initializing the value function in PPO with PRM (**ValueInit**). In the **DenseReward** approach, PRM assigns additional reward signals at each end-of-line token ( $\backslash n$ ) in the code response for coding prompts. Thus, the RL optimization objective for coding prompts is modified to the weighted sum of  $R_{UT}$  and  $R_{PRM}$ , as defined below:

$$\max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}_{train}} \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x})} [R_{UT}(\mathbf{x}, \mathbf{y}) + \lambda R_{PRM}(\mathbf{x}, \mathbf{y}) - \beta \text{KL}(\pi_{\theta} \parallel \pi_{ref})], \quad (3)$$

where  $\lambda$  controls the relative importance of PRM in shaping the reward. Specifically, we set  $\lambda = 0.25$  when the code response does not pass all unit tests, i.e.,  $R_{UT}(\mathbf{x}, \mathbf{y}) = 0$ , and  $\lambda = 0.025$  when the response passes all unit tests, i.e.,  $R_{UT}(\mathbf{x}, \mathbf{y}) = 1$ . The intuition behind this reward shaping is to leverage PRM to provide informative signals when the RL policy fails to generate a valid solution, while minimizing the risk of PRM over-optimization (Rafailov et al., 2024; Skalse et al., 2022) once a correct solution is found. Our empirical results indicate that this reward shaping strategy performs effectively in our experimental setting. In the **ValueInit** setting, PRM is simply used to initialize the value function in PPO. Notably, these two approaches—**DenseReward** and **ValueInit**—are complementary and can be applied concurrently.

## 4.2 KEY ASPECTS FOR INTEGRATING PRM INTO RL TRAINING

While integrating PRM into RL training might seem straightforward, we found that achieving effective results requires careful attention to several critical factors. In this section, we highlight key implementation details essential for the successful application of PRM in RL training.

### 4.2.1 PRM TRAINING: MORE DATA OR BETTER DATA?

Recent research on LLMs highlights that data quality often outweighs quantity (Gunasekar et al., 2023; Li et al., 2023b). We found the same holds true for PRM training data selection. While automated annotation enables large-scale data generation via model sampling, our experiments showed that increasing volume can sometimes degrade PRM performance in RL. Instead, a smaller, well-curated subset provided better supervision and improved outcomes. For instance, when all sampled responses to a prompt consistently pass or fail unit tests, PRM learns little beyond memorization, limiting generalization. We explored various selection and filtering strategies to mitigate this, as detailed in Section 4.3.

### 4.2.2 RL TRAINING: ALLEVIATING PRM HACKING

Reward model hacking (Skalse et al., 2022) is a well-known issue in RLHF training, where the policy learns to exploit flaws in the reward model to achieve high rewards without genuinely improving the quality of response. Similarly, we observed that PRM is also susceptible to such exploitation. Here we discuss two key practical strategies to mitigate the risk of PRM hacking and ensure the reward signals remain aligned with the intended task objectives.

**PRM Reward Length Normalization.** As described in Section 4.1, PRM provides dense rewards by assigning line-level signals at end-of-line tokens. However, directly using PRM predictions,  $R_{\phi}$ , as the reward signal  $R_{PRM}$  in 3 allows exploitation: the policy can generate excessive lines with positive PRM rewards, artificially inflating the total reward. To prevent this, we apply length normalization. Given a prompt  $\mathbf{x}$  and a response  $\mathbf{y}$  with  $T$  lines,  $\mathbf{y} = (\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(T)})$ , the PRM dense reward at line  $m$  is:

$$R_{PRM}(\mathbf{y}^{(m)}) = \frac{1}{T} \cdot R_{\phi}(\mathbf{x}, \mathbf{y}^{\leq m}),$$

which ensures rewards remain bounded in  $[-1, 1]$ , preventing the policy from gaining an advantage by generating excessively long responses.

**Neutral Labeling in PRM Training.** While length normalization curbs reward inflation, models can still exploit PRM by generating excessive comments, which are easier to write than correct code.

To address this, we introduce a neutral label in PRM annotation, as an extension of Eq. (1):

$$l_m = \begin{cases} +1, & \text{if any } \mathcal{P}_k \text{ passes all unit tests} \\ 0, & \text{if the line is a comment} \\ -1, & \text{otherwise} \end{cases}$$

By assigning a neutral label (0) to comments, we eliminate the incentive to generate unnecessary comments, ensuring PRM rewards only meaningful code contributions.

### 4.3 MAIN RESULTS AND ANALYSIS

Table 1: Model performance comparison (Pass@1) on HumanEval, MBPP, and LiveCodeBench datasets. The first section presents the results of public models, alongside the models used in our experiments (Qwen2.5-7B and InHouse-Lite series). The second and third sections detail the performance of RL models under different configurations of PRM usage. Notably, for both models, the best overall performance is achieved when PRM is applied for both Dense Reward and Value Initialization (ValueInit).

Model	Setting		Dataset					
	Dense	Value	HumanEval	MBPP	LiveCodeBench			
	Reward	Init.			Easy	Medium	Hard	Overall
GPT-4o-mini	-	-	87.2	71.8	81.9	27.2	3.6	40.7
Qwen2-72B	-	-	64.6	76.9	65.0	21.3	2.8	32.2
DeepseekCoder-33B	-	-	79.3	70.0	60.8	14.8	1.2	27.7
Qwen2.5-7B-SFT	-	-	67.8	58.1	50.7	16.5	0.9	24.9
InHouse-Lite-SFT	-	-	59.3	59.9	55.3	9.3	0.3	23.5
Qwen2.5-7B-RL	×	×	73.8	62.4	60.9	13.7	1.4	27.5
	×	✓	75.4	63.1	62.8	<b>17.1</b>	<b>1.7</b>	29.6
	✓	×	<b>76.0</b>	63.4	63.1	14.5	1.1	28.5
	✓	✓	74.3	<b>65.4</b>	<b>66.3</b>	15.3	<b>1.7</b>	<b>30.1</b>
InHouse-Lite-RL	×	×	65.1	61.9	<b>70.0</b>	7.2	1.7	28.2
	×	✓	69.8	63.3	67.9	8.9	1.9	28.2
	✓	×	70.0	62.1	68.5	9.9	<b>2.5</b>	28.9
	✓	✓	<b>70.9</b>	<b>63.8</b>	69.3	<b>12.0</b>	1.6	<b>29.8</b>

**Comparing PRM Integration Strategies in RL Training.** We evaluate three strategies for incorporating PRM into RL training, as outlined in Section 4.1: *DenseReward*, *ValueInit*, and a combined approach (*DenseReward & ValueInit*). Table 1 presents the performance of RL models trained using these strategies on HumanEval, MBPP, and LiveCodeBench, alongside SFT and RL baselines. For reference, we also include results from publicly available models such as GPT-4o-mini OpenAI (2023), Qwen2-72B Bai et al. (2023), and DeepseekCoder-33B Guo et al. (2024).

Our results indicate that using PRM as a dense reward signal significantly improves performance over the RL baseline (see the first and third settings for Qwen2.5-7B-RL and InHouse-Lite-RL in Table 1), aligning with findings from Wang et al. (2024a). The granular feedback provided by PRM facilitates policy exploration by offering continuous corrections at intermediate steps. Additionally, using PRM solely for value function initialization also yields consistent improvements. In Qwen2.5-7B-RL, this setting outperforms the RL baseline across all benchmarks, while in InHouse-Lite-RL, similar gains are observed on HumanEval and MBPP.

Combining PRM for both dense rewards and value initialization yields significant performance improvements. In Qwen2.5-7B-RL, Pass@1 increases from 62.4% to 65.4% on MBPP and from 27.5% to 30.1% on LiveCodeBench. Similarly, in InHouse-Lite-RL, Pass@1 improves from 65.1% to 70.9% on HumanEval, 61.9% to 63.8% on MBPP, and 28.2% to 29.8% on LiveCodeBench. This improvement stems from the complementary roles of PRM: dense rewards facilitate exploration by providing rich intermediate feedback, while value initialization stabilizes training and enhances credit assignment. Together, these mechanisms accelerate convergence toward optimal solutions, driving the observed performance gains.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

**PRM Encourages Exploration and Improves Learning Efficiency.** We evaluated the Best-of-K performance for all four training configurations in InHouse-Lite-RL experiments on the training set. Specifically, we assessed all RL models using a decoding configuration with a temperature of 1.0, nucleus sampling (top- $p = 0.95$ ), and top-k sampling ( $k = 128$ ). For each  $K$ , we recorded the percentage of problems solved within  $K$  generated responses, referred to as the *Pass Rate*. Figure 3 presents the results for  $K$  ranging from 1 to 30. Both *DenseReward* and *ValueInit* individually enhance Best-of-K performance compared to the baseline. When combined, they yield the highest improvement, with a nearly 4% increase in Pass Rate at  $K = 30$  over the baseline, highlighting the synergy between dense rewards and value initialization.

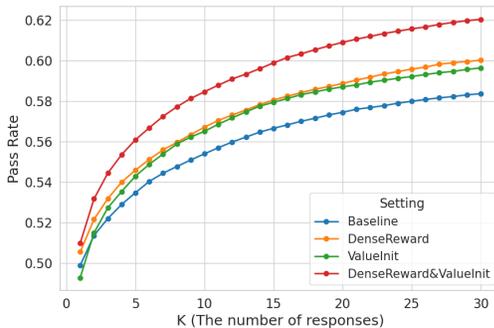


Figure 3: Best-of-K performance curves for all RL training settings, showing the percentage of problems solved within  $K$  generated responses.

**PRM Enhances Code Generation in Long-Horizon Scenarios.** To understand when PRM provides the greatest benefit, we analyze its impact based on response length. Intuitively, the dense nature of the reward signals provided by PRM is particularly advantageous for long-horizon problems, where intermediate feedback can guide policy exploration more effectively.

To validate this, we compared Pass@1 performance of models trained with and without PRM across different response lengths, as shown in Figure 4. Overall, PRM-trained models achieve a 9% improvement in Pass@1 over the baseline. Notably, PRM consistently enhances performance for responses exceeding 100 tokens, whereas for shorter responses, its effect is neutral or slightly negative. We hypothesize that in short-horizon problems, PRM behaves similarly to a biased ORM, offering limited improvements since these problems are already well-explored by the policy. In contrast, for complex, long-horizon problems, PRM provides valuable intermediate signals that help the policy navigate the solution space more effectively, achieving better results with the same optimization compute.

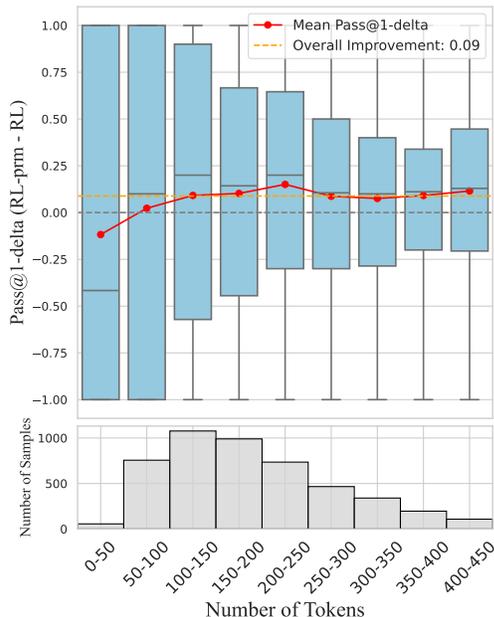


Figure 4: Pass@1 difference between policies trained with and without PRM across varying response lengths.

**The Importance of PRM Training Data Selection** PRM training data can be categorized at two levels: At the *response level*, responses are classified as **Correct** (passes unit tests immediately), **Revised** (initially fails but can find a correct prefix), and **Wrong** (cannot find any correct prefix by binary search within the given budget). At the *prompt level*, prompts are categorized as **Easy** (all responses are **Correct**), **Medium** (mixed response types), and **Hard** (all responses are **Wrong**). We tested the following data selection strategies: **Full** (use all collected data), **Remove Hard** (exclude **Hard** prompts and their responses), **Medium Only** (include only prompts with mixed response types), and **Revised Only** (use only **Revised** responses). We empirically found that **Revised Only**, which includes the richest process-level correction signals, performs best in our setting.

Table 2: Comparison of different PRM data selection strategies on LiveCodeBench (LCB) with InHouse-Lite-RL models.

	Full	Remove Hard	Medium Only	Revised Only
LCB	26.9	27.8	26.9	29.8

**How much data needed to train a PRM that benefits RL training?** Given that automatic PRM data collection is computationally expensive, we examine how the performance of policies trained with PRM scales with the number of training samples. Figure 5 shows how the pass rate of InHouse-Lite-RL models trained with varying amounts of PRM data changes along the average number of responses collected per prompt for PRM data collection, as described in Section 3.1.1. The key finding is that the performance of models trained with PRM improves consistently as the number of PRM training samples increases.

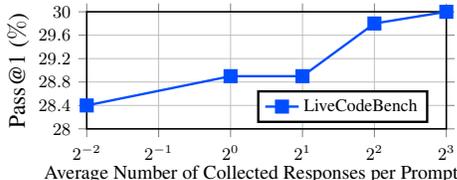


Figure 5: Pass@1 on LiveCodeBench as the average number of responses per prompt for PRM data collection increases (logarithmic scale). A value of  $< 2^0$  indicates that we subsampled prompts from the full dataset.

## 5 RELATED WORKS

### 5.1 LLMs FOR CODE GENERATION

Recently, large language models (LLMs) have demonstrated impressive capabilities in code generation by pre-training on vast text datasets that include code (Lu et al., 2021; Christopoulou et al., 2022; Allal et al., 2023; Zheng et al., 2024; Li et al., 2023b). Additionally, models fine-tuned through supervised fine-tuning (SFT) have achieved competitive results in code generation tasks (Chen et al., 2021a; Li et al., 2023a; Luo et al., 2023; Rozière et al., 2024; Guo et al., 2024). Reinforcement Learning (RL) optimizes policies by interacting with an environment and receiving rewards (Williams, 1992). Recently, RL has been incorporated into LLMs to enhance code generation using unit test feedback (Shojaee et al., 2023; Liu et al., 2023; Le et al., 2022). CodeRL (Le et al., 2022) applies unit test signals as rewards with an actor-critic method, while PPOCoder (Shojaee et al., 2023) builds on this by using the PPO algorithm. RLTF (Liu et al., 2023) improves precision by locating errors, though the reward space remains sparse. Despite progress, RL’s potential to significantly boost code generation in sparse reward environments remains underexplored.

### 5.2 PROCESS REWARD MODELS

Process reward models (PRMs) have garnered significant attention in recent LLM developments, particularly in the mathematical reasoning domain, where they provide verification for intermediate reasoning steps (Lightman et al., 2023; Wang et al., 2024a; Jiao et al., 2024; Wang et al., 2024b; Luo et al., 2024). While some approaches rely on costly and resource-intensive human-annotated process data (Lightman et al., 2023), recent research has focused on automating the collection of process supervision data Wang et al. (2024a); Jiao et al. (2024); Wang et al. (2024b); Luo et al. (2024). Building on these efforts, we similarly automate process supervision but differ in our primary objective. Rather than using PRMs solely as enhanced verifiers compared to Outcome Reward Models (ORMs), we focus on their integration into RL training for code generation. While Wang et al. (2024a) provides preliminary results on PRMs improving RL training in the mathematical domain, their findings are limited. Our work offers a more thorough and systematic investigation of how PRMs can be leveraged in RL for code generation tasks.

## 6 CONCLUSIONS

In this work, we explore a possible way to alleviate the sparse-reward problem in RL-based code generation by introducing a PRM that provides dense, line-level feedback. This approach, inspired by human-like code refinement, enhances learning efficiency and encourages better exploration. Our experiments show that integrating PRMs significantly improves the pass rates of code generation models across HumanEval, MBPP, and LiveCodeBench. Notably, PRM not only facilitates more effective RL training, but also improves performance in long-horizon code generation scenarios.

However, several limitations remain. Its effectiveness depends heavily on data quality, yet automated collection through binary search and unit tests can miss subtle correctness and introduce noise in complex tasks. Training also remains computationally demanding despite binary search optimization. Finally, reliance on external verification (e.g., unit tests) limits applicability to less structured domains like creative writing. Future work should improve data collection and develop alternative evaluation frameworks to extend PRM beyond verifiable domains such as code.

## REFERENCES

- 486  
487  
488 Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz  
489 Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, Logesh Kumar Umapathi,  
490 Carolyn Jane Anderson, Yangtian Zi, Joel Lamy Poirier, Hailey Schoelkopf, Sergey Troshin,  
491 Dmitry Abulkhanov, Manuel Romero, Michael Lappert, Francesco De Toni, Bernardo García  
492 del Río, Qian Liu, Shamik Bose, Urvashi Bhattacharyya, Terry Yue Zhuo, Ian Yu, Paulo Villegas,  
493 Marco Zocca, Sourab Mangrulkar, David Lansky, Huu Nguyen, Danish Contractor, Luis Villa, Jia  
494 Li, Dzmitry Bahdanau, Yacine Jernite, Sean Hughes, Daniel Fried, Arjun Guha, Harm de Vries,  
495 and Leandro von Werra. Santacoder: don't reach for the stars!, 2023. URL <https://arxiv.org/abs/2301.03988>.
- 496  
497 Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David  
498 Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. Pro-  
499 gram synthesis with large language models. *CoRR*, abs/2108.07732, 2021. URL <https://arxiv.org/abs/2108.07732>.
- 500  
501 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
502 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 503  
504 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared  
505 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
506 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
507 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavar-  
508 ian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plap-  
509 pert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol,  
510 Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William  
511 Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan  
512 Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Pe-  
513 ter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech  
514 Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021a.  
URL <https://arxiv.org/abs/2107.03374>.
- 515  
516 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared  
517 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,  
518 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
519 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavar-  
520 ian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plap-  
521 pert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol,  
522 Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William  
523 Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan  
524 Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Pe-  
525 ter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech  
526 Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021b.  
URL <https://arxiv.org/abs/2107.03374>.
- 527  
528 Fenia Christopoulou, Gerasimos Lampouras, Milan Gritta, Guchun Zhang, Yinpeng Guo, Zhongqi  
529 Li, Qi Zhang, Meng Xiao, Bo Shen, Lin Li, Hao Yu, Li Yan, Pingyi Zhou, Xin Wang, Yuchi Ma,  
530 Ignacio Iacobacci, Yasheng Wang, Guangtai Liang, Jiansheng Wei, Xin Jiang, Qianxiang Wang,  
531 and Qun Liu. Pangu-coder: Program synthesis with function-level language modeling, 2022.  
URL <https://arxiv.org/abs/2207.11280>.
- 532  
533 Shihan Dou, Yan Liu, Haoxiang Jia, Limao Xiong, Enyu Zhou, Junjie Shan, Caishuang Huang, Wei  
534 Shen, Xiaoran Fan, Zhiheng Xi, et al. Stepcoder: Improve code generation with reinforcement  
535 learning from compiler feedback. *arXiv preprint arXiv:2402.01391*, 2024.
- 536  
537 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth  
538 Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital  
539 Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai,  
Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need. *CoRR*, abs/2306.11644, 2023. doi: 10.  
48550/ARXIV.2306.11644. URL <https://doi.org/10.48550/arXiv.2306.11644>.

- 540 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
541 Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the  
542 large language model meets programming – the rise of code intelligence, 2024. URL <https://arxiv.org/abs/2401.14196>.  
543
- 544 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
545 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
546 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.  
547
- 548 Fangkai Jiao, Chengwei Qin, Zhengyuan Liu, Nancy F. Chen, and Shafiq Joty. Learning  
549 planning-based reasoning by trajectories collection and process reward synthesizing. *CoRR*,  
550 abs/2402.00658, 2024. doi: 10.48550/ARXIV.2402.00658. URL <https://doi.org/10.48550/arXiv.2402.00658>.  
551
- 552 Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy  
553 Liang. Spoc: Search-based pseudocode to code. In Hanna M. Wallach, Hugo Larochelle,  
554 Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 11883–11894, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/7298332f04ac004a0ca44cc69ecf6f6b-Abstract.html>.  
555
- 556 Hung Le, William La Cava, and Jason H Moore. Coderl: Mastering code generation through pre-  
557 trained models and deep reinforcement learning. *arXiv preprint arXiv:2207.01780*, 2022.  
558
- 559 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao  
560 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii,  
561 Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João  
562 Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Lo-  
563 gesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra  
564 Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey,  
565 Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luc-  
566 cioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor,  
567 Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex  
568 Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva  
569 Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes,  
570 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source  
571 be with you!, 2023a. URL <https://arxiv.org/abs/2305.06161>.  
572
- 573 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.  
574 Textbooks are all you need II: phi-1.5 technical report. *CoRR*, abs/2309.05463, 2023b. doi: 10.  
575 48550/ARXIV.2309.05463. URL <https://doi.org/10.48550/arXiv.2309.05463>.  
576
- 577 Yujia Li, David Choi, Junyoung Chung, Alan Guo, Tsung-Yen Yang, et al. Competition-level code  
578 generation with alphacode. *Science*, 378(6624):1092–1097, 2022.  
579
- 580 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee,  
581 Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *CoRR*,  
582 abs/2305.20050, 2023. doi: 10.48550/ARXIV.2305.20050. URL <https://doi.org/10.48550/arXiv.2305.20050>.  
583
- 584 Jiatae Liu, Yiqin Zhu, Kaiwen Xiao, Qiang Fu, Xiao Han, Wei Yang, and Deheng Ye. RLTF: re-  
585 inforcement learning from unit test feedback. *Trans. Mach. Learn. Res.*, 2023, 2023. URL  
586 <https://openreview.net/forum?id=hjYmsV6nXZ>.  
587
- 588 Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin  
589 Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou,  
590 Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu  
591 Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code understanding  
592 and generation, 2021. URL <https://arxiv.org/abs/2102.04664>.  
593

- 594 Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun  
595 Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. Improve mathematical reasoning in language  
596 models by automated process supervision. *CoRR*, abs/2406.06592, 2024. doi: 10.48550/ARXIV.  
597 2406.06592. URL <https://doi.org/10.48550/arXiv.2406.06592>.
- 598 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
599 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with  
600 evol-instruct, 2023. URL <https://arxiv.org/abs/2306.08568>.
- 601 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774.  
602 URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 603 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
604 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser  
605 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan  
606 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In  
607 *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
608 blfede53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/blfede53be364a73914f58805a001731-Abstract-Conference.html).
- 609 QwenTeam. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.  
610 github.io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 611 Rafael Rafailov, Yaswanth Chittipedu, Ryan Park, Harshit Sikchi, Joey Hejna, Bradley Knox, Chelsea  
612 Finn, and Scott Niekum. Scaling laws for reward model overoptimization in direct alignment  
613 algorithms. *arXiv preprint arXiv:2406.02900*, 2024.
- 614 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
615 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Ev-  
616 timov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong,  
617 Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier,  
618 Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024.  
619 URL <https://arxiv.org/abs/2308.12950>.
- 620 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
621 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL [http://arxiv.org/abs/  
622 1707.06347](http://arxiv.org/abs/1707.06347).
- 623 Parshin Shojaee, Aneesh Jain, Sindhu Tipirneni, and Chandan K Reddy. Execution-based code  
624 generation using deep reinforcement learning. *arXiv preprint arXiv:2301.13816*, 2023.
- 625 Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and character-  
626 izing reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471, 2022.
- 627 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhi-  
628 fang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations.  
629 In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual  
630 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024,  
631 Bangkok, Thailand, August 11-16, 2024, pp. 9426–9439. Association for Computational Lin-  
632 guistics, 2024a. doi: 10.18653/V1/2024.ACL-LONG.510. URL [https://doi.org/10.  
633 18653/v1/2024.acl-long.510](https://doi.org/10.18653/v1/2024.acl-long.510).
- 634 Zihan Wang, Yunxuan Li, Yuexin Wu, Liangchen Luo, Le Hou, Hongkun Yu, and Jingbo Shang.  
635 Multi-step problem solving through a verifier: An empirical analysis on model-induced process  
636 supervision. *CoRR*, abs/2402.02658, 2024b. doi: 10.48550/ARXIV.2402.02658. URL <https://doi.org/10.48550/arXiv.2402.02658>.
- 637 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
638 learning. *Machine learning*, 8:229–256, 1992.
- 639 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
640 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
641 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
642 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng

648 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
649 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
650 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
651 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
652 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

653 Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen,  
654 Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. Codegeex: A pre-trained model  
655 for code generation with multilingual benchmarking on humaneval-x, 2024. URL [https://](https://arxiv.org/abs/2303.17568)  
656 [arxiv.org/abs/2303.17568](https://arxiv.org/abs/2303.17568).  
657

658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A THE USE OF LARGE LANGUAGE MODELS

Large Language Models (LLMs) were used only to improve the manuscript’s presentation by refining language, grammar, readability, and flow. They were not involved in ideation, methodology, experiments, or analysis; all scientific content was developed by the authors. The authors take full responsibility for the manuscript and confirm that LLM use was limited to linguistic polishing in accordance with ethical standards, without contributing to plagiarism or misconduct.

## B A TYPICAL EXAMPLE OF THE LEARNED LINE-WISE REWARDS

In Figure 6, we present a typical example of the line-wise rewards identified by binary search and predicted by a learned PRM to give readers a clearer understanding of our method. In this example, we first sampled a problem from the training set and used our in-house model to generate a response for it. For this generated response (which is not included in the PRM training data), we show the line-wise rewards derived from two sources:

- Line-wise Rewards Identified by Binary Search:** We directly applied the model to perform Algorithm 1, labeling the reward for each line.
- Line-wise Rewards Predicted by a Learned PRM:** We used the learned PRM to predict the rewards for each line.



Figure 6: Visualization of the learned line-wise rewards. The top gray block displays the problem description, while the bottom section shows a model-generated response with line-wise rewards from different sources. The bottom-left block presents the line-wise rewards identified by binary search, and the bottom-right block presents the line-wise rewards predicted by a learned PRM. The actual reward value is shown at the beginning of each line, and each line is color-coded based on the reward value: lines with rewards closer to -1 are shaded red, while those closer to +1 are shaded green.

## C RL TRAINING CURVES

In Figure 7, we present the smoothed RL training curves for all four settings (with and without DenseReward, and with and without ValueInit) in `InHouse-Lite-RL` experiments, using a moving average to reduce noise and enhance readability. These curves correspond to all four RL settings reported in Table 1. The smoothed trends clearly show that when PRM is used as DenseReward, the model solves more problems compared to the baseline, demonstrating PRM’s role in enabling more efficient exploration during RL training. Furthermore, when PRM is applied as both DenseReward and ValueInit, our method achieves the best performance.

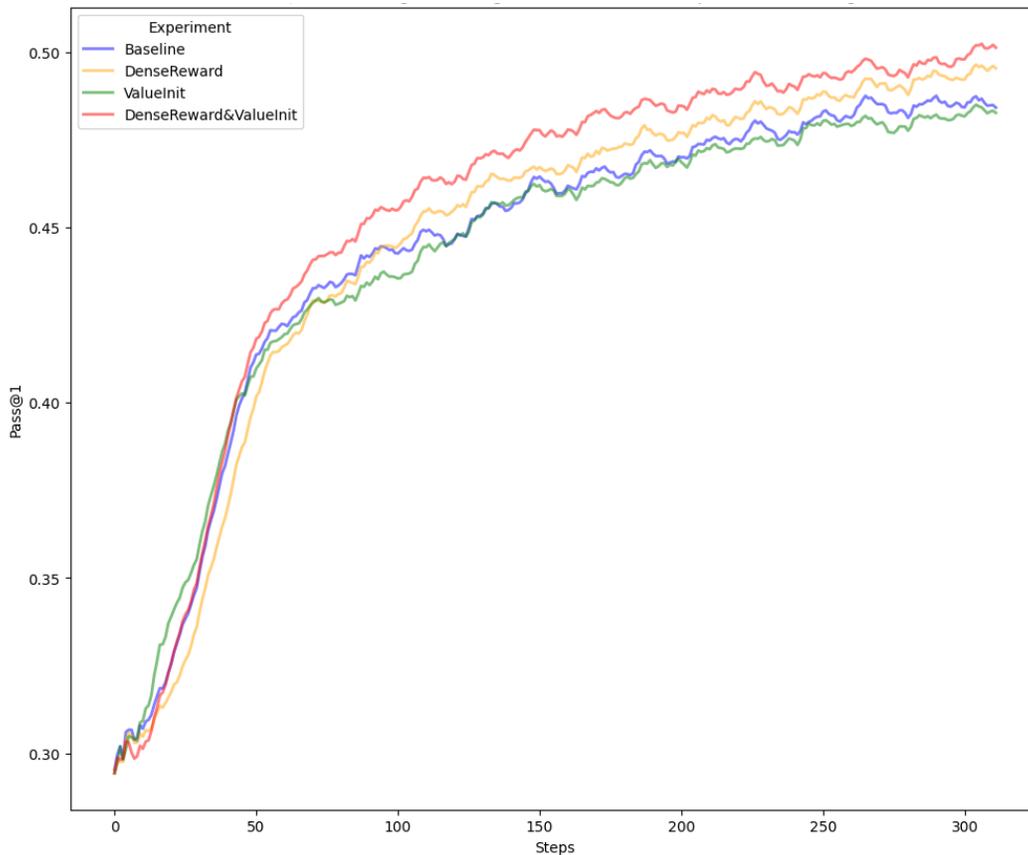


Figure 7: RL training curve of all experiment settings in `InHouse-Lite-RL` experiments. Using PRM as both DenseReward and ValueInit (DenseReward&ValueInit) yields the best result.

## D PRM TRAINING DATA STATISTICS

We present detailed statistics on the PRM datasets used in our experiments to evaluate the optimal PRM data selection strategy, as discussed in Section 4.3. The experiments were conducted using the `InHouse-Lite` series models. Table 3 summarizes the following key metrics: the number of prompt-response pairs (**#Samples**); the total number of tokens across all responses (**#Tokens**); the average number of lines in all responses (**Avg. #Lines**); and the distribution of PRM labels (-1/0/+1).

In Figure 8, we present the distribution of error positions of all **Revised** responses (responses that initially fail but have a correct prefix identified) as determined by the Binary Search procedure (Algorithm 1). The absolute error position (i.e., the position of the first token rejected by Binary Search) is normalized as follows: for a response  $y = (y_1, y_2, \dots, y_L)$  with  $L$  tokens, if the Binary Search accepted the prefix  $(y_1, y_2, \dots, y_p)$  consisting of  $p$  tokens, the Relative Error Position is calculated as  $\frac{p}{L}$ .

Table 3: Statistics of PRM training data collected using different data selection strategies.

Strategy	#Samples	#Tokens	Avg. #Lines	PRM Labels		
				-1	0	+1
<b>Full</b>	838K	179M	16.82	44.25%	17.87%	37.88%
<b>Remove Hard</b>	630K	119M	15.06	24.03%	19.18%	56.79%
<b>Medium Only</b>	485K	104M	16.57	27.66%	19.41%	52.93%
<b>Revised Only</b>	352K	76M	16.71	13.20%	19.42%	67.38%

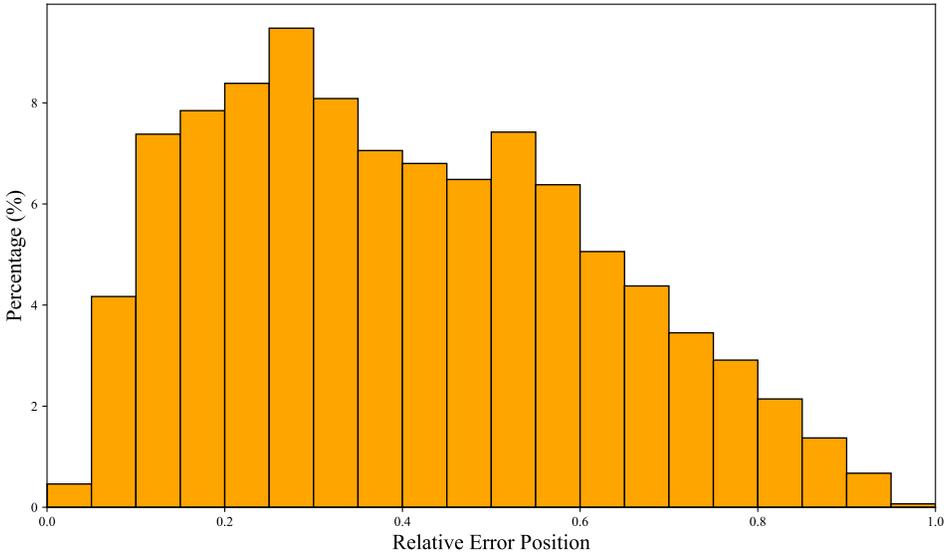


Figure 8: Distribution of Relative Error Positions Identified by Binary Search.

## E QWEN2.5-7B EXPERIMENT DETAILS

**Base Model.** We adopt **Qwen2.5-7B** as our base model (QwenTeam, 2024), a recently released causal language model available at <https://huggingface.co/Qwen/Qwen2.5-7B>. Qwen2.5 belongs to the Qwen series of large language models (Yang et al., 2024), known for their advanced capabilities across a wide range of domains. The Qwen2.5-7B model has 7.61 billion total parameters (6.53 billion excluding embeddings) and utilizes the Transformer architecture as its core. It incorporates state-of-the-art enhancements, including Rotary Positional Embedding (RoPE), SwiGLU activation, RMSNorm, and Attention QKV bias. The model consists of 28 layers and employs 28 attention heads for queries (Q) and 4 for keys and values (KV), making it highly efficient for tasks requiring robust attention mechanisms.

**SFT Settings.** We fine-tuned the Qwen2.5-7B model on the SFT dataset described in Section 4.1. The model was trained for two epochs, starting with a learning rate of  $1 \times 10^{-7}$ , which linearly increased to  $2 \times 10^{-5}$  during the first 2% of the total training steps. After reaching the peak learning rate, a cosine learning rate decay schedule was applied, gradually reducing the learning rate to  $2 \times 10^{-6}$  for the remainder of the training. Additionally, a constant weight decay of 0.01 was used throughout the SFT training process to regularize the model and improve generalization. The model fine-tuned through this process is referred to as **Qwen2.5-7B-SFT**.

**RL Baseline.** We adopted the same RL baseline training method and used the same RLHF dataset described in Section 4.1 to further train the Qwen2.5-7B-SFT model. For PPO training, we configured the following hyperparameters: a batch size of 4096, a linear warmup over the first 5 steps, followed by a constant learning rate of  $2 \times 10^{-6}$  for both the actor and critic, and a KL penalty of 0.01. The training utilized the AdamW optimizer and spanned approximately 300 steps, during which we empirically observed performance convergence.

**PRM Training.** We selected four checkpoints at 50, 100, 150, and 200 steps during the training process of the RL baseline model. For each checkpoint, we sampled  $n = 5$  responses for every coding prompt in the training dataset  $\mathcal{D}_{\text{train}}$ . Each sampled response was labeled using the binary search procedure described in Algorithm 1, with  $K = 20$  completions generated for each partial code prefix. The data collected from all checkpoints was then aggregated to form a PRM training set, employing the **Revised Only** strategy described in Section 4.3. This resulted in 165K samples and 28M tokens. On average, each response contained 16.07 lines. The PRM label distribution was 25.88% for -1, 15.90% for 0, and 58.22% for +1. The PRM was initialized using the value model from the RL baseline and fine-tuned on this PRM dataset using the objective function defined in (2).

**Integrating PRM into RL.** We used the same settings and hyperparameters as described in Section 4.1. Additionally, we observed that due to the properties of the Qwen2.5-7B tokenizer, a newline token is not always represented as a simple "\n" token. Instead, the tokenizer combines other non-space characters with an ending "\n" to form new tokens (e.g., ":\n", "):\n", ")\n", "\n\n", "()\n", "]\n", "()\n", "():\n", etc.). This makes it more challenging to accurately identify line separator tokens in the model’s responses.

Empirically, we addressed this challenge by selecting the 50 most frequent tokens in the PRM dataset whose corresponding token strings include "\n". The full list of token ids is shown below:

{198, 510, 982, 340, 271, 2398, 921, 741, 3932, 1171, 692, 1305, 4167, 2546, 1447, 10343, 1138, 19324, 341, 5563, 9957, 382, 3407, 3646, 624, 48443, 280, 456, 2533, 3989, 1248, 5613, 8389, 8997, 698, 24135, 317, 7368, 2440, 10907, 22165, 4432, 5929, 7129, 345, 11043, 532, 4660, 21686, 14288}.

During RL training, we only applied partial rewards from PRM to these tokens.