

LEARNING TO GENERATE UNIT TEST VIA ADVERSARIAL REINFORCEMENT LEARNING

Dongjun Lee^{1,†}, Changho Hwang², Kimin Lee¹
dgjun32@kaist.ac.kr

KAIST¹, Microsoft Research²

ABSTRACT

Unit testing is a core practice in programming, enabling systematic evaluation of programs produced by human developers or large language models (LLMs). Given the challenges in writing comprehensive unit tests, LLMs have been employed to automate unit test generation, yet methods for training LLMs to produce high-quality unit tests remain underexplored. In this work, we propose UTRL, a novel reinforcement learning (RL) framework that trains an LLM to generate high-quality unit test given a programming instruction. Our key idea is to iteratively train two LLMs, the unit test generator and the code generator, in an adversarial manner via RL: (1) the unit test generator is trained to maximize a discrimination reward, encouraging it to produce tests that reveal faults in the code generator’s solutions; and (2) the code generator is trained to maximize a code reward, encouraging it to produce solutions that pass the unit tests generated by the unit test generator. In our experiment, we demonstrate that unit tests generated by Qwen3-4B trained via UTRL show higher quality compared to unit tests generated by the same model trained via supervised fine-tuning on ground-truth unit tests, yielding code evaluations that more closely align with those induced by the ground-truth tests. Moreover, Qwen3-4B trained with UTRL outperforms frontier models like GPT-4.1 and GPT-4o in generating high-quality unit tests, highlighting the effectiveness of UTRL in training LLMs for the unit test generation.

1 INTRODUCTION

Unit test is a critical component in programming, as it enables evaluation and verification on the functional correctness of the program produced by human developers or large language models (LLMs). Recently, unit tests are widely used as verifiable reward functions in reinforcement learning (RL) or test-time scaling for LLM-driven code generation, where the unit tests provide task-specific feedback signals for guiding the generation of functionally correct code (Le et al., 2022; Guo et al., 2024).

However, implementing unit tests for each programming task is labor-intensive and challenging, since (1) a unit test should contain functionally valid test cases, and (2) each test case in the unit test should cover challenging edge cases, capable of discriminating subtly faulty code implementations. These require significant level of code reasoning capability and understanding of the programming task.

As LLMs have shown strong code understanding and generation capabilities, recent works have explored training them to generate high-quality unit tests given a programming task. A common ap-

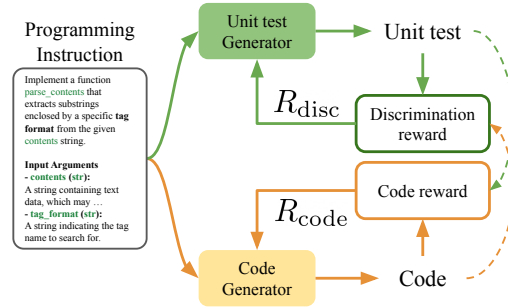


Figure 1: Overview of UTRL. The unit test generator is trained to generate unit test that detect fault in code generated by the code generator, and the code generator is trained to produce code that passes the generated unit test.

[†] Work done during internship at Microsoft Research.

proach is to collect instruction–unit test pairs, followed by supervised fine-tuning (SFT). To enable scalable annotation of high-quality unit tests, these works propose strategies such as using capable teacher models (Ma et al., 2025), applying code perturbation techniques (Prasad et al., 2025), or adopting generator–validator frameworks (Wang et al., 2025b). While promising, SFT-based methods are difficult to scale across diverse programming domains as they fundamentally require unit test labels, which are costly to obtain, for every training example. In contrast, RL removes the need for explicit unit test labels by directly optimizing models against reward signals, offering a more scalable paradigm for training LLMs in unit test generation. The key challenge, however, lies in designing a reward function that can reliably assess the quality of generated unit tests without relying on ground-truth annotations.

In this work, we present **UTRL**, an adversarial RL framework iterating over training a unit test generator LLM and a code generator LLM in an adversarial manner, each LLM defining a reward signal for its counterpart. Our main contribution lies in designing the reward for training unit test generation (i.e., discrimination reward), which rewards the unit test generator LLM when the generated unit tests successfully discriminate the code solution produced by the code generator LLM from ground-truth code solution. As a result, UTRL eliminates the need for ground-truth unit test annotations, since the reward is defined by ground-truth code solution and code generated by the code generator LLM. Specifically, as illustrated in Figure 1, the unit test generator is optimized to maximize the discrimination reward, which encourages producing unit tests that can detect code generated by the code generator LLM from the ground-truth code, while the code generator is optimized to generate code that passes all test cases in the generated unit tests. Through this adversarial training, the code generator LLM progressively learns to generate code solutions that are harder to distinguish from the ground-truth code solutions, and the unit test generator LLM, in turn, learns to generate highly discriminative test cases that can detect subtle faults in the near-correct code solutions, thereby becoming capable of covering challenging edge cases.

In our experiments, we evaluate the quality of generated unit tests on the TACO evaluation set (Li et al., 2023), which contains competitive programming tasks collected from multiple online judge platforms. We first show that unit tests generated by Qwen3-4B (Yang et al., 2025) trained with UTRL, when used as a code reward function for best-of-N sampling, induce $3.1\times$ higher code accuracy gain, compared to the unit tests generated by the base Qwen3-4B. Additionally, we demonstrate that Qwen3-4B trained via UTRL achieves unit test quality higher than the same model trained via SFT, despite not relying on unit test annotations from humans or more capable models, and even surpasses frontier models such as GPT-4.1 (OpenAI, 2025). Finally, the code generator adversarially trained with the unit test generator achieves code generation performance comparable to the same model trained to maximize the pass rate over ground-truth unit tests.

2 RELATED WORK

Unit test generation with LLM Unit testing has been extensively explored for reliable software development (Fraser & Arcuri, 2011; Alagarsamy et al., 2024). As LLMs show strong capability in programming, recent works have explored the automation of unit test generation using large language models (LLMs), with a focus on evaluating the test generation capability of LLMs (Jain et al., 2024; Wang et al., 2024), iteratively refining the unit test by leveraging coverage analysis or mutation testing as a contextual feedback (Alshahwan et al., 2024; Altmayer Pizzorno & Berger, 2025; Dakhel et al., 2024; Foster et al., 2025), or utilizing co-evolutionary algorithm (Li et al., 2025). Another line of work investigates training LLMs for unit test generation via supervised learning, with a focus on constructing high-quality unit test annotations at scale. CodeRM (Ma et al., 2025) leverages stronger teacher models to generate unit tests automatically, while UTGEN (Prasad et al., 2025) proposes a data collection strategy that perturbs code solutions to create faulty variants, and then collects high-quality test cases that discriminate between the correct and faulty code. The work most closely related to ours is the concurrent study CURE (Wang et al., 2025a), which explores the co-evolution of code generation and test case generation capabilities via reinforcement learning. However, our approach differs from CURE in two key aspects: (1) Unlike CURE, which assumes access to a dataset annotated with ground-truth unit tests, UTRL requires only instruction-code pair dataset, which is readily available at large scale (Li et al., 2023; Lambert et al., 2024; Prasad et al., 2025), and (2) UTRL is a framework aimed at training LLMs to generate an optimal set of test cases, rather than a single test case.

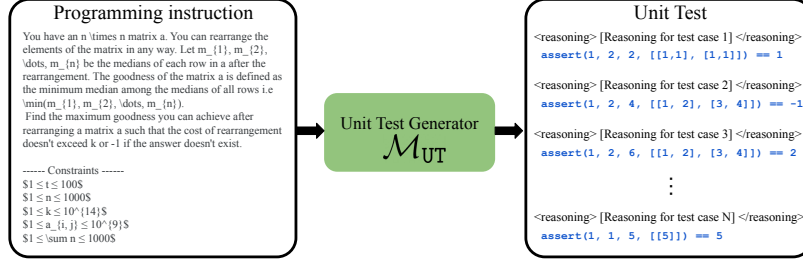


Figure 2: Illustration of input and desired output of unit test generator LLM. Given a programming instruction specifying input arguments and functionality of a code, the unit test generator LLM generates a set of N test cases comprehensively covering the various edge cases based on reasoning.

Reinforcement learning for LLM Early works on training LLMs with reinforcement learning (RL) utilized reward models learned from human preference feedback and trained LLM with the Proximal policy optimization algorithm (PPO) (Schulman et al., 2017) for aligning LLMs with human values, establishing reinforcement learning with human feedback (RLHF) as a standard framework for instruction tuning and alignment in LLMs (Stiennon et al., 2020; Ouyang et al., 2022). To address the complexity and instability of online RL, subsequent methods proposed offline preference learning algorithms that bypass the need for reward models while maintaining strong alignment performance (Rafailov et al., 2023; Ethayarajh et al., 2024; Hong et al., 2024). More recently, RL algorithms with improved efficiency have been proposed as a variant of PPO, and such RL algorithms have been actively applied to improve the reasoning capabilities of LLMs, particularly in domains such as math and programming where verifiable rewards can be precisely defined (Shao et al., 2024; Guo et al., 2024; Yu et al., 2025). However, different from code generation or math, defining the verifiable rewards for unit test generation is non-trivial and remains underexplored.

3 METHOD

In this section, we present UTRL, a framework for adversarially training a unit test generator LLM and a code generator LLM via RL. Section 3.1 describes an overview of the UTRL, and the core components of UTRL are described in Section 3.2 and Section 3.3.

3.1 OVERVIEW

Our goal is to train a unit test generator LLM \mathcal{M}_{UT} , which takes a programming instruction I as input and outputs a unit test \mathcal{T} , as shown in Figure 2. A unit test \mathcal{T} consists of multiple test cases, i.e., $\mathcal{T} = \{T_i\}$, where each T_i denotes a single test case. Each test case T_i checks the partial correctness of a code implementation C for instruction I , by verifying that C produces the expected output for a given input. A straightforward approach to training \mathcal{M}_{UT} is to collect a dataset of instruction-unit test pairs and fine-tune LLMs via supervised learning. However, curating such a dataset is inherently challenging because implementing unit tests is an open-ended task with no unique solution, and it is often difficult to objectively evaluate their quality.

To address these challenges, we propose a framework for training a unit test generator LLM \mathcal{M}_{UT} without relying on instruction-unit test pairs. Since instruction-code solution pairs are widely available, we train the \mathcal{M}_{UT} via RL to produce unit tests that can distinguish ground-truth code from imperfect code generated by a code generator LLM, \mathcal{M}_{code} . We further extend this idea into an adversarial framework: \mathcal{M}_{UT} improves by learning to generate unit test identifying weaknesses in code generated by \mathcal{M}_{code} , while \mathcal{M}_{code} in turn improves by learning to produce code that passes increasingly challenging unit tests.

Specifically, we iterate over the following two key steps (see Algorithm 1):

- **Step 1. Training the unit test generator \mathcal{M}_{UT} :** For a given programming instruction I , we first sample multiple code solutions using the code generator LLM \mathcal{M}_{code} . The unit test generator LLM \mathcal{M}_{UT} is then trained to produce unit tests that (1) discriminate as many sampled code solutions from the ground-truth solution as possible, and (2) contain functionally valid test cases (see Section 3.2).

Algorithm 1 UTRL

Require: Dataset of instruction-code solution pairs $\mathcal{D} = \{(I_k, C_k^*)\}$,
Initial code generator LLM $\mathcal{M}_{\text{code}}$, Initial unit test generator LLM \mathcal{M}_{UT} ,
Weight for the discrimination reward λ ,
Desired minimum number of test cases per unit test τ

```

1: while not converged do
2:   // Training unit test generator  $\mathcal{M}_{\text{UT}}$  (see Section 3.2)
3:   for  $(I, C^*) \in \mathcal{D}$  do
4:      $\mathcal{T} \sim \mathcal{M}_{\text{UT}}(\cdot | I)$  ▷ Generate unit test  $\mathcal{T}$  using  $\mathcal{M}_{\text{UT}}$ 
5:      $\mathcal{C} \sim \mathcal{M}_{\text{code}}(\cdot | I)$  ▷ Generate a set of code solutions  $\mathcal{C}$  using  $\mathcal{M}_{\text{code}}$ 
6:      $r_{\text{UT}} = \lambda R_{\text{disc}}(\mathcal{T}, \mathcal{C}, C^*) + (1 - \lambda) R_{\text{valid}}(\mathcal{T}, C^*, \tau)$  ▷ Compute reward for  $\mathcal{T}$  (see Equation 1, 2)
7:      $\mathcal{M}_{\text{UT}} := \text{RL-update}(\mathcal{M}_{\text{UT}}, r_{\text{UT}})$  ▷ Update the unit test generator  $\mathcal{M}_{\text{UT}}$ 
8:   end for
9:   // Training code generator  $\mathcal{M}_{\text{code}}$  (see Section 3.3)
10:  for  $(I, C^*) \in \mathcal{D}$  do
11:     $C \sim \mathcal{M}_{\text{code}}(\cdot | I)$  ▷ Generate code  $C$  using  $\mathcal{M}_{\text{code}}$ 
12:     $\mathcal{T} \sim \mathcal{M}_{\text{UT}}(\cdot | I)$  ▷ Generate unit test  $\mathcal{T}$  using  $\mathcal{M}_{\text{UT}}$ 
13:     $r_{\text{code}} = R_{\text{code}}(C, \mathcal{T}, C^*)$  ▷ Compute reward for  $C$  (see Equation 4)
14:     $\mathcal{M}_{\text{code}} := \text{RL-update}(\mathcal{M}_{\text{code}}, r_{\text{code}})$  ▷ Update code generator  $\mathcal{M}_{\text{code}}$ 
15:  end for
16: end while

```

- **Step 2. Training code generator $\mathcal{M}_{\text{code}}$:** Given the programming instruction I , the code generator $\mathcal{M}_{\text{code}}$ is trained to generate solutions that maximize the pass rate over \mathcal{T} , which is generated by \mathcal{M}_{UT} (see Section 3.3).

By repeating these steps, the code generator progressively learns to produce more functionally correct solutions, while the unit test generator learns to generate increasingly discriminative and high-quality unit tests.

3.2 TRAINING UNIT TEST GENERATOR

To train the unit test generator LLM \mathcal{M}_{UT} to produce both effective and functionally valid unit tests, we introduce two reward components. The discrimination reward R_{disc} serves as the primary objective, encouraging \mathcal{M}_{UT} to generate unit tests that can detect faults across diverse incorrect code solutions. To simultaneously ensure the functional validity of each test case, we introduce the validity reward R_{valid} , which measures the proportion of functionally valid test cases among the entire test cases in the generated unit test. During training, the unit test generator \mathcal{M}_{UT} is optimized with the weighted sum of R_{disc} and R_{valid} via RL.

Discrimination reward The discrimination reward evaluates how effectively the test cases in \mathcal{T} discriminate LLM-generated code solutions from ground-truth code solutions. Given a unit test \mathcal{T} generated regarding programming instruction I , the discrimination reward corresponding to \mathcal{T} is formally defined as follows:

$$R_{\text{disc}}(\mathcal{T}, \mathcal{C}, C^*) = \frac{1}{|\mathcal{C}|} \sum_{C \in \mathcal{C}} \left[1 - \prod_{T \in \mathcal{T}} (1 - \text{Pass}(C, T))^{\text{Pass}(C^*, T)} \right], \quad (1)$$

where the C^* is a ground-truth code solution, \mathcal{C} is a set of code solutions generated by the code generator $\mathcal{M}_{\text{code}}$, and $\text{Pass}(C, T)$ is an indicator function that returns 1 if the code C passes the test case T , and 0 otherwise.⁰ To compute the discrimination reward regarding unit test \mathcal{T} , as illustrated in Figure 3, we first filter out functionally invalid test cases from the generated unit test \mathcal{T} , ensuring that invalid test cases do not affect the discrimination reward. We then leverage a set of code solutions \mathcal{C} generated by $\mathcal{M}_{\text{code}}$ (i.e., $\mathcal{C} \sim \mathcal{M}_{\text{code}}(\cdot | I)$), where the discrimination reward is computed as a fraction of the code solutions that do not pass the filtered unit test.

Validity reward The validity reward evaluates whether each test case in the generated unit test \mathcal{T} is functionally valid (i.e., correctly maps the input and expected output of the code solution), and is formally defined as follows:

$$R_{\text{valid}}(\mathcal{T}, C^*, \tau) = \frac{\sum_{T \in \mathcal{T}} \text{Pass}(C^*, T)}{\max(|\mathcal{T}|, \tau)}, \quad (2)$$

⁰See Appendix A.3 for details about the implementation of Pass function.

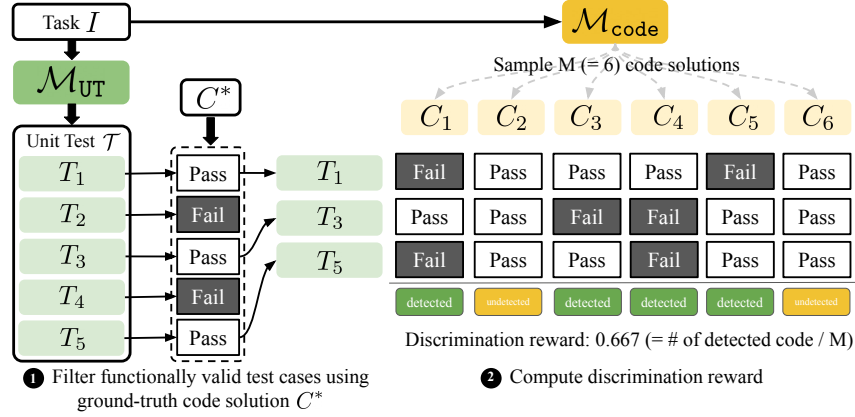


Figure 3: Overview of the process of computing discrimination reward with respect to unit test \mathcal{T} . First, among the 5 test cases in the unit test \mathcal{T} , test cases that pass under the ground-truth code C^* (i.e., T_1, T_3, T_5) are filtered, forming a set of functionally valid test cases. Second, the discrimination reward is defined as a ratio of sampled code solutions that do not pass at least one valid test case. In this figure, among 6 sampled code solutions, 4 code solutions (C_1, C_3, C_4, C_5) do not pass at least one valid test case, resulting in the discrimination reward of 0.667 ($= \frac{4}{6}$).

where τ is a hyperparameter to adjust the desired minimum number of test cases in a single generated unit test. We apply clipping to the denominator in Equation 2 as $\max(|\mathcal{T}|, \tau)$, in order to prevent unit tests with only a small number of trivial test cases from receiving an inappropriately high validity reward. To be specific, if the R_{valid} is simply defined as a fraction of valid test cases among total test cases (i.e., $R_{valid}(\mathcal{T}, C^*) = \sum_{T \in \mathcal{T}} \text{Pass}(C^*, T) / |\mathcal{T}|$), unit tests with an extremely small number of test cases (e.g., unit tests with a single trivial test case) acquire high validity reward, which is undesirable. To mitigate this problem, we clip the denominator as $\max(|\mathcal{T}|, \tau)$. This ensures that unit tests with fewer than τ test cases receive a validity reward proportional to the absolute number of valid test cases, thereby preventing the unit tests with a small number of trivial test cases from receiving high validity reward.

Update unit test generator LLM In order to guide the unit test generator LLM to produce unit test comprising test cases that are both (1) highly discriminative and (2) functionally valid, we define a training reward as a weighted sum of the discrimination reward R_{disc} and the validity reward R_{valid} :

$$r_{UT} = \lambda R_{disc}(\mathcal{T}, C, C^*) + (1 - \lambda) R_{valid}(\mathcal{T}, C^*, \tau), \quad (3)$$

where the λ is a hyperparameter that weights the discrimination reward.

3.3 TRAINING CODE GENERATOR

After training the unit test generator LLM \mathcal{M}_{UT} , we train the code generator LLM \mathcal{M}_{code} via RL. Regarding a code C generated by \mathcal{M}_{code} , the training reward is formally defined as follows:

$$R_{code}(C, \mathcal{T}, C^*) = \frac{\sum_{T \in \mathcal{T}} \text{Pass}(C, T) \cdot \text{Pass}(C^*, T)}{\sum_{T \in \mathcal{T}} \text{Pass}(C^*, T)}, \quad (4)$$

where C^* is a ground-truth code solution. In detail, we first filter out functionally invalid test cases (i.e., test cases that do not pass under the ground-truth code solution C^*) from \mathcal{T} , because such invalid test cases lead to faulty evaluation of the code. We then measure the proportion of the functionally valid test cases that are passed by the generated code C . Based on the reward design, we optimize \mathcal{M}_{code} to produce code solutions that pass unit tests generated by \mathcal{M}_{UT} .

4 EXPERIMENTS

We design our experiments to investigate the following questions:

- **RQ1.** Is UTRL effective at training LLMs to generate high-quality unit tests? (see Section 4.2)
- **RQ2.** Is UTRL more effective than supervised learning-based approaches? (see Section 4.3)

- **RQ3.** Is UTRL effective at improving code generation? (see Section 4.4)
- **RQ4.** Does iterative training of UTRL enable continuous improvement of unit test generation capability? (see Section 4.5)

4.1 EXPERIMENTAL SETUP

In this section, we describe training details, baselines, and evaluation protocols for our experiments.

Training details As the base model for both the code generator and unit test generator, we use Qwen3-4B (Yang et al., 2025). For RL training, we adopt Grouped Relative Policy Optimization (GRPO; Shao et al. 2024 (see Appendix A.1 for details) and use 15,249 programming instruction-code solution pairs from the TACO dataset (Li et al., 2023). Further details about training and prompts are described in Appendix A.3 and A.5.

Baselines We compare our method to various baselines for unit test generation. Implementation details for the baselines are provided in Appendix A.3.

- **LLM baselines without fine-tuning:** We use 4 open-source LLMs (Qwen3-4B, Qwen3-8B, Qwen-14B, Deepseek-Coder-V2-Lite; Qwen et al. 2025; Zhu et al. 2024) and 2 closed LLMs (GPT-4o, GPT-4.1; Hurst et al. 2024; OpenAI 2025) as baselines for unit test generation.
- **Supervised fine-tuning (SFT):** We fine-tune Qwen3-4B with instruction-unit test pairs through supervised learning. For this purpose, we use high-quality unit tests from the TACO training dataset \mathcal{D}_{UT} . To investigate whether reasoning signals can improve unit test generation, we construct a reasoning-augmented dataset $\mathcal{D}_{\text{reason}+UT}$ containing unit tests paired with their corresponding reasoning. Specifically, we first generate both reasoning and unit tests for each training task using Gemini-2.5-flash (Google Deepmind, 2025), capable LLMs in code generation.¹ We then filter out functionally invalid test cases and retain the valid reasoning-unit test pairs as target labels for SFT training.
- **CURE:** We also consider CURE (Wang et al., 2025a), an RL algorithm that fine-tunes an LLM to perform both test case generation and code generation using a dataset of instruction-unit test pairs. Specifically, we evaluate unit tests generated by ReasonFlux-Coder-7B (Wang et al., 2025a), which is based on Qwen2.5-7B-Instruct (Yang et al., 2025) fine-tuned with CURE.
- **Ground-truth (GT) unit tests:** We also compare against the unit tests provided in the TACO evaluation dataset (Li et al., 2023), which we consider as ground truth (GT) and regard as an upper bound since they are rigorously verified and consist of a large number of high-quality test cases.

Evaluation protocols For evaluation, we utilize rigorously verified 945 competitive programming tasks provided in the TACO evaluation dataset (Li et al., 2023), which were collected from various online judge platforms (CodeForces, CodeChef, HackerRank, and HackerEarth). To measure the quality of the generated unit tests, we propose the following evaluation schemes:

- **Best-of- N improvement:** To examine whether the generated unit tests include high-quality test cases, we perform best-of- N sampling for code generation, using the generated unit tests as evaluators. Specifically, we generate N candidate solutions from code LLMs and select the one that passes the largest number of generated unit tests. To evaluate the quality of the selected code, we then measure two metrics: (1) code score, defined as the fraction of test cases in a GT unit test (i.e., high-quality unit tests from the TACO evaluation dataset) that the code solution passes, and (2) code accuracy, which indicates whether the solution passes all test cases in the GT unit test.
- **Unit test fidelity:** To quantify how closely a generated unit test approximates the evaluation induced by the GT unit test, we introduce a metric called *unit test fidelity*. For each task, we first sample multiple code solutions from LLMs. We then evaluate these solutions twice: once using the generated unit test and once using the GT unit test. This yields two code score vectors, and we compute Spearman’s correlation between them. A high correlation indicates that the generated unit test closely replicates the evaluation induced by the comprehensive GT unit test.

¹Instead of attaching reasoning to the original unit tests in TACO, which led to frequent hallucinations and inconsistencies, we regenerate both reasoning and unit tests together.

Unit test generated by	Code LLM: Qwen3-8B (32 shot)				Code LLM: Qwen3-14B (32 shot)			
	Score	Δ Score	Acc (%)	Δ Acc (%p)	Score	Δ Score	Acc (%)	Δ Acc (%p)
Qwen3-4B	0.430	+0.084	9.8	+1.8	0.459	+0.057	13.4	+2.7
Qwen3-8B	0.435	+0.089	9.7	+1.7	0.470	+0.070	12.6	+1.9
Qwen3-14B	0.460	+0.114	9.8	+1.8	0.483	+0.083	11.6	+0.9
Deepseek-Coder-V2-Lite	0.433	+0.087	9.5	+1.5	0.452	+0.050	11.9	+1.2
GPT-4o	0.457	+0.111	10.6	+2.6	0.494	+0.094	12.1	+1.4
GPT-4.1	0.528	+0.182	13.7	+6.2	0.548	+0.148	14.8	+4.6
Qwen3-4B + SFT w/ \mathcal{D}_{UT}	0.458	+0.112	11.7	+3.7	0.491	+0.091	14.0	+3.3
Qwen3-4B + UTRL (ours)	<u>0.558</u>	<u>+0.212</u>	<u>14.9</u>	<u>+6.9</u>	<u>0.588</u>	<u>+0.188</u>	<u>17.0</u>	<u>+6.3</u>
Qwen3-14B + UTRL (ours)	0.562	+0.216	15.7	+7.7	0.594	+0.194	17.1	+6.4
GT (upperbound)	0.650	+0.304	21.0	+13.1	0.680	+0.278	24.4	+13.7

Table 1: Best-of- N improvement achieved by Qwen3-4B and Qwen3-14B trained via UTRL, compared against baselines. We report code score (Score) and code accuracy (Acc) of the best-of- N selected code solution, and also report the increment of each metric compared to code generated without the best-of- N sampling (i.e., Δ Score and Δ Acc). This result demonstrates effectiveness of UTRL in training LLMs to produce unit tests with highly discriminative test cases.

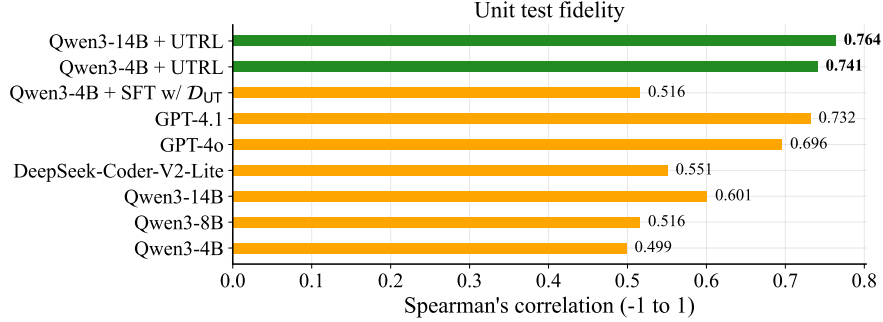


Figure 4: Fidelity of unit tests generated by Qwen3-4B and Qwen3-14B trained with UTRL, compared against baselines. Both model trained via UTRL achieves the unit test fidelity higher than baselines, demonstrating the effectiveness of UTRL in training LLMs to generate unit tests that closely approximates the code evaluation induced by the GT unit tests.

We provide more details about the evaluation in Appendix A.4.

4.2 QUALITY OF UNIT TESTS

In this section, we demonstrate the effectiveness of UTRL in training LLMs to generate high-quality unit tests, compared with strong baselines. We first assess the quality of the generated unit tests using Best-of- N improvement and unit test fidelity. We then compare UTRL against a recent RL-based baseline, CURE.

Best-of- N improvement Table 1 reports the best-of- N improvement achieved by Qwen3-4B and Qwen3-14B trained with UTRL, compared against several baselines. When used as code evaluators for best-of-32 sampling with Qwen3-8B and Qwen3-14B code generators, unit tests produced by Qwen3-4B trained with UTRL yield code accuracies of 14.9% and 17.0%, respectively. For comparison, Qwen3-4B trained with SFT achieves only 11.7% and 14.0%. Moreover, Qwen3-4B trained with UTRL outperforms frontier proprietary models such as GPT-4.1 and GPT-4o, underscoring the effectiveness of UTRL in training LLMs to generate high-quality unit tests. With increased model capacity, Qwen3-14B trained with UTRL attains stronger performance, reaching code accuracies of 15.7% and 17.1%. Notably, UTRL is effective even when paired with closed-source code generators such as GPT-4o (see Table 3 in Appendix A.2 for supporting results).

Unit test fidelity Figure 4 presents the fidelity of unit tests generated by Qwen3-4B trained with UTRL, compared against several baselines. Qwen3-14B and Qwen3-4B trained with UTRL achieves a Spearman’s correlation of 0.764 and 0.741, even outperforming GPT-4.1, demonstrat-

Unit test generated by	Code LLM: Qwen3-8B (32 shot)				Code LLM: Qwen3-14B (32 shot)			
	Score	Δ Score (%)	Acc (%)	Δ Acc (%p)	Score	Δ Score	Acc (%)	Δ Acc (%p)
CURE	0.446	+0.100	9.7	+1.7	0.483	+0.083	12.7	+2.0
UTRL [†] (ours)	0.521	+0.175	14.1	+6.1	0.562	+0.162	15.9	+5.2
GT (upperbound)	0.650	+0.304	21.0	+13.1	0.680	+0.278	24.4	+13.7

Table 2: Best-of- N improvement induced by UTRL and CURE. For UTRL, we train Qwen2.5-7B-Instruct using 4.5K programming tasks in CodeContest dataset (Li et al., 2022) following the training setup of CURE, which we denote as UTRL[†].

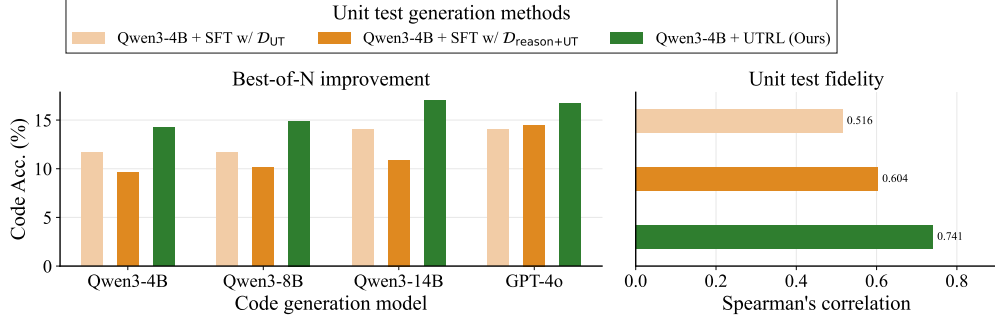


Figure 5: Best-of- N improvmenet (**left**) and unit test fidelity (**right**) achieved by UTRL, compared against SFT baselines (SFT w/ \mathcal{D}_{UT} , SFT w/ $\mathcal{D}_{reason+UT}$).

ing the effectiveness of UTRL in training LLMs to generate unit tests that induce code evaluations consistent with the GT unit tests. Interestingly, while SFT with \mathcal{D}_{UT} provides substantial gains in best-of- N improvement relative to the base Qwen3-4B, it yields only marginal improvements in unit test fidelity. We conjecture that this is because SFT with \mathcal{D}_{UT} trains the LLM to generate unit tests without step-by-step reasoning, making it difficult for the model to produce logically structured test suites (i.e., ranging from basic test cases to highly discriminative edge cases), which are essential for achieving high unit test fidelity.

Comparison to CURE Table 2 and Figure 6 report the best-of- N improvement and unit test fidelity achieved by Qwen2.5-7B-Instruct trained with UTRL versus the same model trained with CURE, an RL baseline based on instruction–unit test pairs. When used as a code evaluator for best-of-32 sampling with Qwen3-8B and Qwen3-14B code generators, UTRL achieves code accuracies of 14.1% and 15.9%, which are 4.4% and 3.2% higher than those obtained with CURE. Moreover, UTRL achieves a unit test fidelity of 0.593, surpassing CURE. These results show that UTRL achieves superior performance despite relying only on instruction–code pairs, in contrast to CURE requiring instruction–unit test pairs.

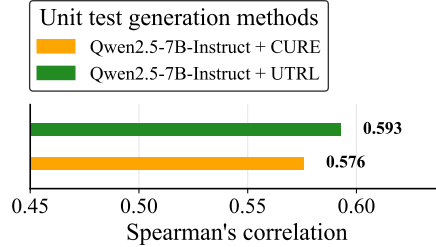


Figure 6: Unit test fidelity of UTRL compared against CURE. Qwen2.5-7B-Instruct trained via UTRL achieves unit test fidelity higher than CURE, demonstrating superiority of UTRL over CURE.

4.3 SUPERIORITY OF UTRL OVER SUPERVISED LEARNING APPROACHES

To investigate whether training LLMs with reasoning-annotated unit tests via SFT can improve unit test generation, we also evaluate SFT with $\mathcal{D}_{reason+UT}$ and compare it with UTRL. Figure 5 presents best-of- N improvement (left) and unit test fidelity (right) for Qwen3-4B trained with SFT baselines versus the same model trained with UTRL. While SFT + $\mathcal{D}_{reason+UT}$ achieves an improved unit test fidelity score of 0.604, higher than the SFT + \mathcal{D}_{UT} , it falls behind 0.741 achieved by UTRL. Although SFT baselines directly optimize LLMs using unit test labels annotated by humans or more capable teacher models, they show limited generalization to evaluation tasks compared with UTRL. This observation is consistent with prior findings that SFT tends to memorize the training distribu-

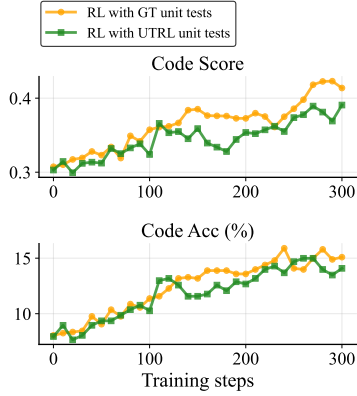


Figure 7: Code score and accuracy on the TACO evaluation tasks. We evaluate Qwen3-4B code generator trained with UTRL (green) and compare it with the same model trained using rewards defined as the pass rate over GT unit tests (orange).

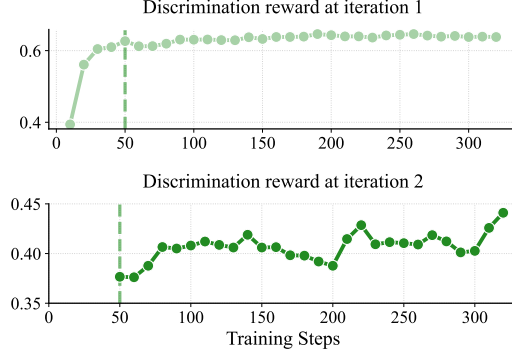


Figure 8: Discrimination reward evaluated at training iterations 1 and 2. During iteration 1 training, the discrimination reward saturated after 50 steps with only about a gain of 0.02, while in iteration 2, it exhibited an improvement of 0.07.

tion, whereas RL promotes better generalization in reasoning tasks (Chu et al., 2025). Moreover, because collecting high-quality unit test annotations for SFT is costly, UTRL offers a more scalable and practical approach for training LLMs for unit test generation.

4.4 EFFECTIVENESS OF UTRL IN TRAINING CODE GENERATOR

In UTRL, the code generator is trained jointly with the unit test generator to produce solutions that maximize the pass rate over the generated unit tests. To evaluate the effectiveness of UTRL for training the code generator, we compare Qwen3-4B trained with UTRL against the same model trained with RL using rewards defined as the pass rate over GT unit tests. Figure 7 shows code accuracy and code score averaged over 945 evaluation tasks across RL training steps. The code generator trained with UTRL achieves a code accuracy of 15.3% and a code score of 0.391, comparable to the 15.9% and 0.423 obtained when training with rewards from GT unit tests. These results indicate that UTRL is effective in training code generators, and that the unit tests it produces are sufficiently reliable to serve as rewards for code generation tasks.

4.5 EFFECT OF ITERATIVE TRAINING

We examine whether iterative training creates increasingly challenging discrimination tasks for the unit test generator, and whether the generator improves by learning from them. Figure 8 shows discrimination rewards averaged over evaluation tasks. In the first iteration, the unit test generator is trained against code produced by the initial code generator, and the reward saturates after about 50 steps. At this point, we stop training and update the code generator via RL to maximize its pass rate against the current unit tests. When training resumes in iteration 2, the reward drops from 0.626 to 0.375, indicating that the updated code generator now produces solutions that are harder to distinguish from ground truth, thereby creating a more challenging discrimination task for the unit test generator. As training continues in iteration 2, the unit test generator improves the reward from 0.375 to 0.447. Correspondingly, as shown in Figure 9, the unit test generator at iteration 2 produces unit tests that yield superior best-of- N performance compared to iteration 1.

5 CONCLUSION

In this work, we present UTRL, an adversarial reinforcement learning framework for training LLM as unit test generator. Our approach alternates between training the unit test generator and a code generator, enabling the unit test generator to produce unit tests that can distinguish between near-correct code solutions and the ground-truth code implementation. Through extensive experiments and analysis, we demonstrate the superiority of UTRL over recent baselines and supervised learning-based methods, even outperforming frontier models like GPT-4.1. We hope this work paves the way for future research on training algorithm for unit test generation.

ETHICS STATEMENT

We introduce UTRL, an RL framework for training LLMs for unit test generation. As unit tests enable systematic and interpretable verification over source code written by LLMs, we believe the improvement in automated unit test generation will contribute to more reliable and safe LLM-based code generation and agentic software engineering.

REPRODUCIBILITY STATEMENT

For the reproducibility of our results, we have provided a detailed description of our experimental setups and prompts in Section A.3, A.4, and A.5. Additionally, to further facilitate reproduction, we will open-source the training dataset, model checkpoint, and source code.

REFERENCES

- Saranya Alagarsamy, Chakkrit Tantithamthavorn, and Aldeida Aleti. A3test: Assertion-augmented automated test case generation. *Information and Software Technology*, 176:107565, 2024.
- Nadia Alshahwan, Jubin Chheda, Anastasia Finogenova, Beliz Gokkaya, Mark Harman, Inna Harper, Alexandru Marginean, Shubho Sengupta, and Eddy Wang. Automated unit test improvement using large language models at meta. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*, pp. 185–196, 2024.
- Juan Altmayer Pizzorno and Emery D Berger. Coverup: Effective high coverage test generation for python. *Proceedings of the ACM on Software Engineering*, 2(FSE):2897–2919, 2025.
- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025.
- Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Majdinasab, Foutse Khomh, and Michel C Desmarais. Effective test generation using pre-trained large language models and mutation testing. *Information and Software Technology*, 171:107468, 2024.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Christopher Foster, Abhishek Gulati, Mark Harman, Inna Harper, Ke Mao, Jillian Ritchey, Hervé Robert, and Shubho Sengupta. Mutation-guided llm-based test generation at meta. *arXiv preprint arXiv:2501.12862*, 2025.
- Gordon Fraser and Andrea Arcuri. Evosuite: automatic test suite generation for object-oriented software. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*, pp. 416–419, 2011.
- Google Deepmind. Gemini 2.5 flash. <https://deepmind.google/models/gemini/flash/>, 2025. Accessed: 2025-06-17.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*, 2024.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Kush Jain, Gabriel Synnaeve, and Baptiste Roziere. Testgeneval: A real world unit test generation and test completion benchmark. *arXiv preprint arXiv:2410.00752*, 2024.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.

- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:21314–21328, 2022.
- Kefan Li, Hongyue Yu, Tingyu Guo, Shijie Cao, and Yuan Yuan. Cocodevo: Co-evolution of programs and test cases to enhance code generation. *arXiv preprint arXiv:2502.10802*, 2025.
- Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. Taco: Topics in algorithmic code generation dataset. *arXiv preprint arXiv:2312.14852*, 2023.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022. doi: 10.1126/science.abq1158. URL <https://www.science.org/doi/abs/10.1126/science.abq1158>.
- Zeyao Ma, Xiaokang Zhang, Jing Zhang, Jifan Yu, Sijia Luo, and Jie Tang. Dynamic scaling of unit tests for code reward modeling. *arXiv preprint arXiv:2501.01054*, 2025.
- OpenAI. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>, 2025. Accessed: 2025-04-14.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Archiki Prasad, Elias Stengel-Eskin, Justin Chih-Yao Chen, Zaid Khan, and Mohit Bansal. Learning to generate unit tests for automated debugging. *arXiv preprint arXiv:2502.01619*, 2025.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020.
- Wenhan Wang, Chenyuan Yang, Zhijie Wang, Yuheng Huang, Zhaoyang Chu, Da Song, Lingming Zhang, An Ran Chen, and Lei Ma. Testeval: Benchmarking large language models for test case generation. *arXiv preprint arXiv:2406.04531*, 2024.
- Yinjie Wang, Ling Yang, Ye Tian, Ke Shen, and Mengdi Wang. Co-evolving llm coder and unit tester via reinforcement learning. *arXiv preprint arXiv:2506.03136*, 2025a.

Zihan Wang, Siyao Liu, Yang Sun, Hongyan Li, and Kai Shen. Codecontests+: High-quality test case generation for competitive programming. *arXiv preprint arXiv:2506.05817*, 2025b.

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

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

Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*, 2024.

A APPENDIX

A.1 GROUPED RELATIVE POLICY OPTIMIZATION

In this section, we describe the RL algorithm we used for training the LLMs in our experiments.

Proximal Policy Optimization Proximal Policy Optimization (PPO) (Schulman et al., 2017), an actor-critic RL algorithm, is widely used for training LLMs via RL. PPO updates the LLM policy π_θ by maximizing the clipped surrogate objective described as follows:

$$\mathcal{J}_{\text{PPO}}(\theta) = \mathbb{E}_{o \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\sum_{t=1}^{|o|} \min \left(\frac{\pi_\theta(o_t|x, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t|x, o_{<t})} A(x, o_t), \right. \right. \\ \left. \left. \text{clip} \left(\frac{\pi_\theta(o_t|x, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t|x, o_{<t})}, 1 - \epsilon, 1 + \epsilon \right) A(x, o_t) \right) \right], \quad (5)$$

where π_θ is a policy under the training, x denotes a prompt, $o_{<t}$ is a completion up to t -th token, o_t is a t -th token generated by the LLM $\pi_{\theta_{\text{old}}}$, and ϵ is a hyperparameter for clipping the importance sampling ratio. Here, the computation of advantage $A(x, o_t)$ requires learning a separate value function V_ψ . Usually, such a value function is parameterized by LLM of size comparable to the LLM policy under training, thereby requiring significant memory and computational cost.

Grouped Relative Policy Optimization In order to bypass the cost of learning the value function, Grouped Relative Policy Optimization (GRPO) has been proposed as a variant of PPO (Shao et al., 2024). In order to compute the advantage, GRPO samples a group of G outputs $\{o^{(i)}\}_{i=1}^G$ from the old LLM policy $\pi_{\theta_{\text{old}}}$ for input prompt x , and computes scalar rewards $\{r_i = R(x, o^{(i)})\}_{i=1}^G$. These rewards are then normalized by subtracting the group mean and dividing by the group standard deviation. The resulting normalized score $A_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)}$ is broadcast to all tokens t in the output $o^{(i)}$, such that $A(x, o_t^{(i)}) = A_i$. This shared advantage is then used to weight the token-level policy gradients, allowing for effective optimization without a learned value function. Formally, GRPO optimizes the following objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o^{(i)}|} \sum_{t=1}^{|o^{(i)}|} \min \left(\frac{\pi_\theta(o_t^{(i)}|x, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)}|x, o_{<t}^{(i)})} A_i, \right. \right. \\ \left. \left. \text{clip} \left(\frac{\pi_\theta(o_t^{(i)}|x, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)}|x, o_{<t}^{(i)})}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) \right], \quad (6)$$

where $A_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)}$ is the normalized advantage shared across all tokens in output $o^{(i)}$.

A.2 ADDITIONAL EXPERIMENT

Ablation on validity reward In order to demonstrate the effectiveness of the validity reward and its design choice (see Equation 2), we conduct an ablation study comparing the following variants of the UTRL, (1) UTRL trained without validity reward (i.e., w/o R_{valid}), and (2) UTRL with validity reward but without clipped normalization in validity reward term, where the validity reward is defined as a ratio of valid test cases among the generated test cases (i.e., w/o clipping). Figure 10 shows the number of test cases and the ratio of valid test cases in the generated unit test over training steps. UTRL without R_{valid} results in the generation of unit tests with more than 50% of invalid test cases, while UTRL results in unit tests containing less than 35% of invalid test cases. Furthermore, in the case of UTRL without clipping, we observe that the model collapses to generate an extremely small number of trivial test cases in order to maximize the ratio of the valid test cases, hindering the learning progress toward generating comprehensive unit tests. These results demonstrates the effectiveness of the design choice of UTRL.

Best-of-N improvement measured with additional code generators Table 3 presents Best-of-N improvement measured when Qwen3-4B and GPT-4o is used for code generator LLM. In this evaluation, UTRL still achieves the highest code accuracy and code score, demonstrating that LLMs trained via UTRL produces high-quality unit tests that can effectively evaluate code generated by various code LLMs.

Unit test generated by	Qwen3-4B				GPT-4o			
	Score	Δ Score	Acc (%)	Δ Acc (%p)	Score	Δ Score	Acc (%)	Δ Acc (%p)
Qwen3-4B	0.430	+0.084	9.8	+1.8	0.459	+0.057	13.4	+2.7
Qwen3-8B	0.435	+0.089	9.7	+1.7	0.470	+0.070	12.6	+1.9
Qwen3-14B	0.461	+0.115	9.7	+1.7	0.475	+0.075	13.3	+2.6
DeepSeek-Coder-Lite-V2	0.407	+0.061	9.0	+1.0	0.458	+0.056	12.9	+2.2
GPT-4o	0.457	+0.111	10.6	+2.6	0.494	+0.094	12.1	+1.4
GPT-4.1	0.528	+0.182	14.2	+6.2	0.548	+0.148	15.3	+4.6
Qwen3-4B + SFT w/ \mathcal{D}_{UT}	0.457	+0.127	11.7	+4.5	0.437	+0.239	14.0	+7.5
Qwen2.5-7B + CURE	0.421	+0.092	10.2	+3.0	0.459	+0.261	13.8	+7.3
Qwen2.5-7B + UTRL [†]	0.502	+0.173	14.0	+6.8	0.526	+0.328	16.0	+9.5
Qwen3-4B + UTRL	0.536	+0.207	14.2	+7.0	0.551	+0.353	16.7	+10.2
Qwen3-14B + UTRL	0.548	+0.219	15.0	+7.8	0.557	+0.359	17.5	+11.0
GT (upperbound)	0.650	+0.304	21.0	+13.1	0.680	+0.278	24.4	+13.7

Table 3: Best-of-N improvement measured using Qwen3-4B and GPT-4o as code generator LLM. UTRL shows the highest best-of-N improvement compared to baselines. Here, UTRL [†] denotes the variant of UTRL using 4.5K programming tasks for training, for fair comparison to CURE.

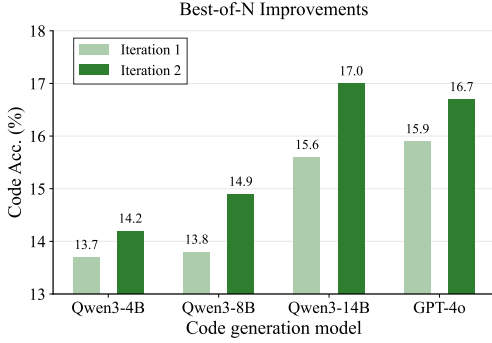


Figure 9: Best-of-N performance gains obtained by UTRL at the same iterations. Across 4 code generation models, iteration 2 results in improved best-of-N improvement compared to iteration 1.

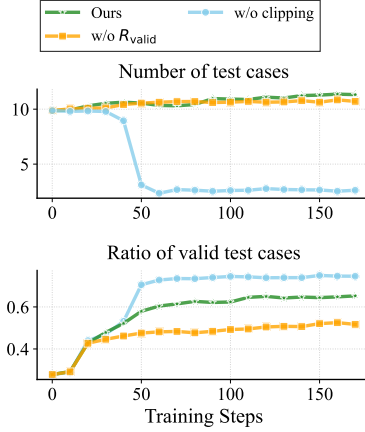


Figure 10: Averaged number of test cases in the generated unit test and ratio of valid test cases across RL training steps. Ablated variants of UTRL results in generation of unit test with invalid test cases, or only few trivial test cases.

A.3 TRAINING DETAILS

In this section, we provide implementation details of UTRL and other baselines in our experiments. Hyperparameters are described in Table 4.

Training details of UTRL For training, we utilize 15,249 programming task instances provided in training split of TACO dataset, where we exclude instances without a ground-truth code solution and those whose unit test annotation do not follow `stdio` format (e.g., functional assertions) or without ground-truth code solution annotation, and 1,000 task instances for validation. We train Qwen3-4B and Qwen3-14B via UTRL. For Qwen3-4B experiment, we train the unit test generator for 50 steps in the first iteration, and train the code generator for 370 steps, and then continually train the unit test generator for 100 steps in the second iteration. For Qwen3-14B experiment, we train the unit test generator for 100 steps at the first iteration, without further iteration. We use weighted sum of discrimination reward and validity reward (see Equation 3) measured on the 1,000 validation tasks as a validation metric to choose training steps. For the number of LLM-generated code samples to define a discrimination reward, we set the value as 8. For τ , a hyperparameter to adjust the desired



Figure 11: Examples of unit test annotations for two datasets for SFT baselines (SFT with \mathcal{D}_{UT} and SFT with $\mathcal{D}_{\text{reason}+UT}$). $\mathcal{D}_{\text{reason}+UT}$ is annotated by ground-truth unit tests, while $\mathcal{D}_{\text{reason}+UT}$ is annotated by the teacher model (Gemini-2.5-flash).

minimum number of the test cases in the unit test, we set the value as 12. Additionally, for λ , a weighting factor for the discrimination reward, we set this value as 0.85.

Implementation of Pass function We describe implementation details about the `Pass` function, which is core component of reward computation. Given a code C and a test case T , we consider $\text{Pass}(C, T)$ is 1 (i.e., code C passes test case T) if (1) the code C builds successfully without any syntax error, (2) the execution of test case T regarding the code C does not return any error (e.g., assertion error), and (3) the execution of test case T regarding the code C runs under 10 seconds. Otherwise, we consider $\text{Pass}(C, T)$ is 0 (i.e., code C does not pass test case T).

Training details of baselines

- **SFT with \mathcal{D}_{UT} :** As a training dataset, we use the 15,249 training task instances in TACO dataset, exactly same as UTRL, and we fine-tune Qwen3-4B. To label the unit test for SFT, we randomly select 12 test cases from the ground-truth unit test for each task in the TACO dataset. As shown in Figure 11 (left), we format the target response by listing the 12 test cases without reasoning and use them directly as SFT labels.
- **SFT with $\mathcal{D}_{\text{reason}+UT}$:** As a training dataset, we use the 15,249 training task instances in TACO dataset, and we fine-tune Qwen3-4B. In order to label $\mathcal{D}_{\text{reason}+UT}$ with reasoning-annotated unit tests, we prompt Gemini-2.5-flash with the programming instruction to generate 12 test cases with reasoning. We then filter out functionally invalid test cases (i.e., test cases that fail under the ground-truth solution), and use the remaining valid test cases annotated with reasoning as SFT labels. We remark that this is replication of CodeRM (Ma et al., 2025), which utilized Llama3-70B-Instruct to label unit test with reasoning, and train the Llama3-8B-Instruct with the labeled dataset.
- **CURE:** For evaluating CURE and compare it with UTRL, we use the open-source ReasonFlux-7B-Coder, which is Qwen2.5-7B-Instruct fine-tuned with the CURE algorithm using 4.5K programming tasks from a subset of the CodeContest dataset (Li et al., 2022). In order to ensure fair comparison, we train Qwen2.5-7B-Instruct with UTRL using the same 4.5K training tasks (denoted as UTRL[†]), and compare the resulting model with ReasonFlux-7B-Coder. Additionally, following the evaluation setup of CURE (Wang et al., 2025a), for each programming instruction in evaluation set, we sample 16 test cases to form a single unit test, using the prompt format provided in the CURE paper.

A.4 EVALUATION DETAILS

Evaluation tasks Across all experiments, we use 945 competitive programming tasks in the evaluation split of the TACO dataset. The evaluation tasks span diverse difficulty levels, and each task is annotated with highly comprehensive unit test (i.e., ground-truth unit test) which comprises an

Methods	Optimizer	LR	Batch size	warmup	kl coef	sampling temperature
UTRL (iter 1)	AdamW	1e-6	128	1e-2	1e-3	1.0
UTRL (iter 2)	AdamW	5e-6	128	1e-2	1e-3	1.0
CURE	AdamW	1e-6	128	1e-2	1e-2	1.0
SFT w/ \mathcal{D}_{UT}	AdamW	1e-5	128	1e-2	-	-
SFT w/ $\mathcal{D}_{reason+UT}$	AdamW	1e-5	128	1e-2	-	-

Table 4: Training hyperparameters of UTRL and baselines.

average of 202.3 test cases. The tasks were collected from Codeforces, CodeChef, HackerRank, and HackerEarth, which are widely used competitive programming platforms, where the task distribution is described in Table 5.

	CodeForces	CodeChef	HackerRank	HackerEarth
# of tasks	540	245	20	40

Table 5: Distribution of evaluation tasks over multiple competitive programming platforms.

Best-of-N improvement For measuring Best-of-N improvement, we utilize LLM for sampling 32 code solutions, and select the best code solution by using the generated unit test as a code evaluator (i.e., select the code solution that passes the largest number of test cases in the generated unit test). We then evaluate code score and code accuracy regarding the selected code solution, using ground-truth unit tests provided in the TACO evaluation dataset. An intuition behind this evaluation is that a more discriminative unit test will identify and select higher-quality code among multiple candidate code solutions. We evaluate the Best-of-N improvement on multiple code generator LLMs, Qwen3-4B, 8B, 14B, and GPT-4o.

Unit test fidelity Unit test fidelity measures whether the generated unit tests induce code evaluation consistent with the code evaluation by the ground-truth unit tests. For each programming task, we first sample 128 code solutions using Qwen3-4B, 8B, 14B and GPT-4o (i.e., sample 32 code solutions using each LLM). We then evaluate code score of the 128 code solutions twice, once by using the ground-truth unit test, and once by using the generated unit test, which induces two score vectors length of 128 (i.e., number of code solutions under evaluation) and each score ranges from 0 to 1. Finally, we compute Spearman’s correlation between these two vectors to quantify how closely the generated unit tests reproduce the code evaluations of the ground-truth unit test.

A.5 PROMPTS

In this section, we provide the prompts used in our experiments: (1) unit test generation, (2) code generation, and (3) rationalization, which were used for the RT baseline in Section 4.3.

Unit test generation prompt

System prompt:

You are an expert Python programmer capable of generating test cases for Python programming tasks.
 Given a programming task, generate several independent test cases with corresponding reasoning.
 Each test case should be independent of the others and sharply target distinct corner cases so that arbitrary faulty code solutions can be detected.
 Before writing each test case, think deeply about the input arguments that expose extreme or subtle edge cases, and reason about the expected output.
 After completing the reasoning process, generate the test case in stdio format.
 Specifically, your output should follow the format below:

```
<reasoning>
Reasoning for test case 1
```

```

First, reason about the input arguments that can discriminate an
incorrect code solution (e.g., edge cases).
Ensure the input arguments test aspects independent of
previous test cases.
Then, derive the expected output from the problem description.
</reasoning>
'''

Input:
stdio format input 1

Output:
stdio format output 1
'''

<reasoning>
Reasoning for test case 2
First, reason about the input arguments that can discriminate an
incorrect code solution (e.g., edge cases).
Ensure the input arguments test aspects independent of
previous test cases.
Then, derive the expected output from the problem description.
</reasoning>
'''

Input:
stdio format input 2

Output:
stdio format output 2
'''

...

<reasoning>
Reasoning for test case 12
First, reason about the input arguments that can discriminate an
incorrect code solution (e.g., edge cases).
Ensure the input arguments test aspects independent of
previous test cases.
Then, derive the expected output from the problem description.
</reasoning>
'''

Input:
stdio format input 12

Output:
stdio format output 12
'''

Ensure the following:
1. Do not include solution code in your response; generate
exactly 12 test cases.
2. For each test case, provide a detailed rationale.
3. Each test case must be independent; avoid duplicates.

User prompt:
Here is the problem description:

{problem_query}

Based on comprehensive reasoning, generate comprehensive unit test
involving several test cases for the given problem.
The test cases should cover various edge cases,
corner cases, and normal cases, at the same time, functionally correct.

```

Code generation prompt**System prompt:**

You are an expert Python programmer.
 Based on the problem description, solve the coding problem efficiently.
 Think step by step, then write a Python solution that solves the problem.
 Follow the format below:

```
<reasoning>
Write your reasoning here.
</reasoning>

```python
Write your Python solution here. It should run with stdio-format input.
```
```

User prompt:

Here is the problem description:

```
{problem_query}
```

A.6 LIMITATIONS & FUTURE DIRECTIONS

We outline the limitations of UTRL and promising directions for future research.

Performance gap with ground-truth unit tests Although we show that UTRL is a promising direction for improving unit test generation capabilities of LLMs, a performance gap remains between UTRL-generated unit tests and ground-truth unit tests. As UTRL can be combined with arbitrary online RL algorithms, we believe that investigation on RL algorithm enabling better exploration, and combination of UTRL with the better RL algorithm will further improve the unit test generation performance.

Extension to broad software engineering domains While our work conducts experiments on the competitive programming task domain, UTRL can be applied to broader software / programming domains. Scaling UTRL with large-scale instruction-code datasets covering a wider range of programming scenarios might be a promising direction for future work.

Considerations for variable length unit tests UTRL trains the LLM to generate fixed number of test cases per unit test. Future work could explore framework that enables adaptive generation of variable-length unit tests depending on the complexity of the given programming task.