EVOSCALE: EVOLUTIONARY TEST-TIME SCALING FOR SOFTWARE ENGINEERING

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

033

034

037

038

040

041 042

043

044

046 047

048

051

052

ABSTRACT

Language models (LMs) perform well on standardized coding benchmarks but struggle with real-world software engineering tasks such as resolving GitHub issues in SWE-Bench—especially when model parameters are less than 100B. While smaller models are preferable in practice due to their lower computational cost, improving their performance remains challenging. Existing approaches primarily rely on supervised fine-tuning (SFT) with high-quality data, which is expensive to curate at scale. An alternative is test-time scaling: generating multiple outputs, scoring them using a verifier, and selecting the best one. Although effective, this strategy often requires excessive sampling and costly scoring, limiting its practical application. We propose Evolutionary Test-Time Scaling (EvoScale), a sampleefficient method that treats generation as an evolutionary process. By iteratively refining outputs via selection and mutation, EvoScale shifts the output distribution toward higher-scoring regions, reducing the number of samples needed to find correct solutions. To reduce the overhead from repeatedly sampling and selection, we train the model to *self-evolve* using reinforcement learning (RL). Rather than relying on external verifiers at inference time, the model learns to self-improve the scores of its own generations across iterations. Evaluated on SWE-Bench-Verified, EvoScale enables a 32B model to match or exceed the performance of models with over 100B parameters while using a few samples. Code, data, and models will be fully open-sourced.

1 Introduction

Language models (LMs) perform well on coding benchmarks like HumanEval (Chen et al., 2021) or LiveCodeBench (Jain et al., 2025a) but struggle with real-world software engineering (SWE) tasks (Jimenez et al., 2024). Unlike standardized coding problems, real issues—such as GitHub issues (Jimenez et al., 2024)—are often under-specified and require reasoning across multiple files. Even large models like Claude reach only around 60% accuracy on SWE-bench (Jimenez et al., 2024), despite using carefully engineered prompting pipelines (Xia et al., 2024). Smaller models (under 100B parameters) perform significantly worse, typically scoring below 10% in zero-shot settings and plateauing around 30% after supervised fine-tuning (SFT) (Xie et al., 2025; Pan et al., 2025) on GitHub issue datasets. Improving the performance of these models remains a key challenge for practical deployment, where repeatedly querying large models is often too costly or inefficient.

Recent and concurrent works to improve the performance of small LMs on SWE tasks have mainly focused on expanding SFT datasets—either through expert annotation or distillation from larger models (Yang et al., 2025; Xie et al., 2025; Pan et al., 2025). These approaches show that performance improves as the quality and quantity of training data increase. However, collecting such data is both costly and time-consuming.

An alternative is *test-time scaling*, which improves performance by generating multiple outputs at inference and selecting the best one using a verifier, such as a reward model (Cobbe et al., 2021; Lightman et al., 2023). While widely applied in math and logical reasoning (Hoffmann et al., 2022; Snell et al., 2025), test-time scaling remains underexplored in SWE. Yet it shows strong potential: prior works (Pan et al., 2025; Brown et al., 2024) demonstrate that small models can generate correct solutions when sampled many times. Specifically, their pass@N, the probability that at least one of N samples is correct, is close to the pass@1 performance of larger models. This indicates that small models can produce correct solutions; the challenge lies in efficiently identifying them.

Test-time scaling assumes that among many sampled outputs, at least one will be correct. However, when correct solutions are rare, these methods often require a large number of samples to succeed. This is particularly costly in SWE tasks, where generating each sample is slow due to long code contexts, and scoring is expensive when unit tests execution is needed (Xia et al., 2024). Recent work (Wei et al., 2025) uses reinforcement learning (RL) (Wei et al., 2025) to enhance the reasoning capabilities of LMs for improved output quality but still requires hundreds of code edits (i.e., patch samples) per issue. Also, Pan et al. (2025) depends on slow interactions with the runtime environment in agentic workflows. This motivates the need for *sample-efficient* test-time scaling methods that can identify correct solutions with fewer samples.

In this paper, we propose Evolutionary Test-Time Scaling (EvoScale), a sample-efficient method for improving test-time performance on SWE tasks. Existing test-time scaling methods often require an excessive number of samples because model outputs are highly dispersed—correct solutions exist but are rare, as shown in Figure 1. EvoScale mitigates this by progressively steering generation toward higher-scoring regions, reducing the number of samples needed to find correct outputs. Inspired by evolutionary algorithms (Shen et al., 2023; Wierstra et al., 2014; Hansen, 2016; Salimans et al., 2017), EvoScale iteratively refines candidate patches through selection and mutation. Instead of consuming the sample budget in a single pass, EvoScale amortizes it over multiple iterations: the model generates a batch of outputs, a scoring function selects the top ones, and the next batch is generated by conditioning on these—effectively mutating prior outputs. Early iterations focuses on exploration; later ones focus on exploitation.

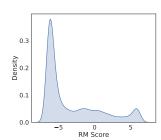


Figure 1: Reward score distribution of SFT outputs, with high-scoring outputs concentrated in the long tail.

Although EvoScale improves sample efficiency, the selection step still incurs overhead: like standard evolutionary algorithms (Wierstra et al., 2014), it generates more outputs than needed and filters only the high-scoring ones, increasing sampling and computation costs. To eliminate this, we use RL to internalize the reward model's guidance into the model itself, enabling it to *self-evolve*—refining its own outputs without external reward models at inference. We formulate this as a potential-based reward maximization problem (Ng et al., 1999), where the model learns to improve output scores across iterations based on score differences. This avoids discarding low-scoring outputs and reduces sample usage per iteration. Our theoretical analysis shows that this RL objective ensures monotonic score improvement across iterations. We evaluate the proposed EvoScale method on SWE-Bench-Verified (Jimenez et al., 2024), and summarize our key contributions as follows:

- A new perspective of formulating test-time scaling as an evolutionary process, improving sample efficiency for software engineering tasks.
- A novel RL training approach that enables self-evolution, eliminating the need for external reward models or verifiers at inference time.
- Empirical results showing that a 32B model with EvoScale achieves performance comparable to models exceeding 100B parameters, while requiring only a small number of samples

2 RELATED WORK

Dataset Curation for SWE. Prior works (Ma et al., 2024; Pan et al., 2025) and concurrent efforts (Yang et al., 2025; Jain et al., 2025b) use proprietary LLMs (e.g., Claude, GPT-4) as autonomous agents to collect SFT data by recording step-by-step interactions in sandboxed runtime environments. While this automates the data collection process for agent-style training (Yang et al., 2024), it involves substantial engineering overhead (e.g., Docker setup, sandboxing) and high inference costs. In contrast, Xie et al. (2025) uses a pipeline-based framework (Xia et al., 2024), collecting real pull-request–issue pairs and prompting GPT-4 to generate CoT traces and ground-truth patches without runtime interaction. Though easier to collect, this data requires careful noise filtering. Our approach instead improves small models' performance by scaling the computation at test time.

Test-time scaling for SWE. Xia et al. (2024) showed that sampling multiple patches and selecting the best one based on unit test results in sandboxed environments improves performance. Unit tests have since been widely adopted in SWE tasks (Wei et al., 2025; Ehrlich et al., 2025; Jain et al., 2025; Brown et al., 2024). Other works (Pan et al., 2025; Jain et al., 2025b; Ma et al., 2025) train

verifiers or reward models to score and select patches. To reduce the cost of interacting with runtime environments in agentic frameworks (Yang et al., 2024), some methods (Ma et al., 2025; Antoniades et al., 2025) integrate tree search, pruning unpromising interaction paths early. While prior works improve patch ranking or interaction efficiency, our focus is on reducing the number of samples needed for effective test-time scaling.

RL for SWE. Pan et al. (2025) used a basic RL approach for SWE tasks, applying rejection sampling to fine-tune models on their own successful trajectories. Wei et al. (2025) later used policy gradient RL (Shao et al., 2024), with rewards based on string similarity to ground truth patches, showing gains over SFT. In contrast, our method trains the model to iteratively refine its past outputs, improving scores over time. We also use a learned reward model that classifies patches as correct or incorrect, which outperforms string similarity as shown in Appendix A.

3 PRELIMINARIES



Figure 2: **Pipeline for SWE Tasks.** Given a GitHub issue, the retriever identifies the code files most relevant to the issue. The code editor then generates a code patch to resolve it.

Software engineering (SWE) tasks. We study the problem of using LMs to resolve real-world GitHub issues, where each issue consists of a textual description and a corresponding code repository. Since issues are not self-contained, solving them requires identifying and modifying relevant parts of the codebase. There are two main paradigms for solving SWE tasks with LMs: agentic (Yang et al., 2024) and pipeline-based (Xia et al., 2024; Wei et al., 2025). Agentic methods allow the model to interact with the runtime environment, such as browsing files, running shell commands, and editing code through tool use. While flexible, these approaches are computationally intensive and rely on long-context reasoning, making them less practical for small models. In contrast, pipeline-based methods decompose the task into subtasks, typically retrieval and editing, and solve each without runtime interaction, which is more computationally efficient and suited for small models. Retrieval refers to identifying the files or functions relevant to the issue, while editing involves generating the code changes needed to resolve it.

Formally, given an issue description x, the goal is to produce a code edit (i.e., patch) y that fixes the bug or implements requested changes. A retrieval model selects a subset code context $C(x) \subseteq \mathcal{C}$ from the full codebase \mathcal{C} , and an editing model π generates the patch $y = \pi(x, C(x))$ by modifying C(x). While retrieval has reached around 70% accuracy in prior work Xie et al. (2025); Xia et al. (2024), editing remains the main bottleneck. This work focuses on improving editing performance in pipeline-based settings, using off-the-shelf localization methods in experiments. In this setup, the dominant cost comes from sampling and scoring outputs from the editing model at test time.

Test-time scaling (Hoffmann et al., 2022; Snell et al., 2025) improves model performance during inference without training. It typically involves sampling multiple outputs and selecting the best one using a scoring function (e.g., a reward model) (Cobbe et al., 2021; Lightman et al., 2023). Specifically, the model generates outputs y_1, \ldots, y_N , scores them with R, and returns $\arg\max_{y_i} R(y_i)$. This strategy is commonly used in domains like reasoning and mathematics.

4 METHOD: EVOLUTIONARY TEST-TIME SCALING

Goal: Sample-efficient test-time scaling. Test-time scaling improves performance by selecting the best output from multiple samples, but often requires a large number of generations to find correct solutions, especially in SWE tasks (Wei et al., 2025). Our goal is to enable test-time scaling more sample-efficient, achieving stronger performance with fewer samples.

Why is test-time scaling sample-inefficient in SWE task? Correct solutions exist but are rarely sampled, as for hard issues, the model's output distribution is not concentrated around high-scoring regions. Given a sample budget N, typical test-time scaling methods in SWE (Xia et al., 2024; Wei et al., 2025; Jain et al., 2025b; Pan et al., 2025) draw N outputs (patches) $\{y_i\}_{i=1}^N$ from a frozen

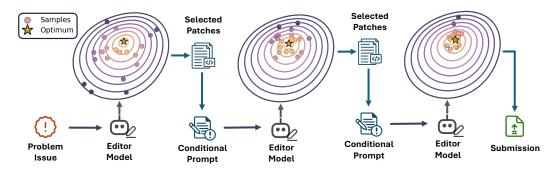


Figure 3: An Overview of Evolutionary Test-Time Scaling. Given a GitHub issue x and its code context C(x), the editor model π first generates a batch of candidate patches \mathcal{Y}^t . The reward landscape is illustrated with contour lines, where brighter contours indicate a higher score of a scoring function R (e.g., reward model or unit tests). A set of patches \mathcal{E}^t is selected (e.g., via a scoring function R) and combined with x and C(x) to form a conditional prompt (see Section 4.1), which guides the model to generate the next batch $\mathcal{Y}^{t+1} = \{y_1^{t+1}, \dots, y_M^{t+1}\}$, increasingly concentrated around the optimum. The process continues under a fixed sampling budget until convergence, after which the final patch is submitted.

editor model π , score them with a score function R (e.g., reward model or unit tests), and selects the best one $\arg\max_{y_i} R(x,y_i)$. While high-scoring outputs near the mode could be sampled easily, the challenge of test-time scaling is to identify high-scoring outputs from the tail of $\pi(\cdot \mid x, C(x))$. However, doing so typically requires a large sample size N, making the process sample-inefficient.

Our approach: This motivates our method, Evolutionary Test-Time Scaling (EvoScale), which iteratively refines generation by using earlier outputs to guide subsequent sampling. We recast patch generation for a GitHub issue as an evolutionary process. The objective is to explore the patch space with a small number of samples, identify high-scoring patches, and iteratively refine the generated patches. As shown in Figure 3, initial samples are scattered and far from the correct solutions (denoted by stars), but over iterations, the distribution shifts closer to the correct solution. Through evolution, EvoScale more efficiently uncovers high-scoring outputs in long tails. We formulate the problem in Section 4.1 and detail the training procedure in Sections 4.2 and 4.3.

4.1 FORMULATION: PATCH GENERATION AS EVOLUTION

We amortize the sampling budget over T iterations by generating M < N samples per iteration, rather than sampling all N at once. The goal is to progressively improve sample quality across iterations. A key challenge lies in effectively using early samples to guide later ones. Typical evolutionary strategies select top-scoring candidates and mutate them—often by adding random noise—to steer future samples toward high-scoring regions. However, in SWE tasks, where patches are structured code edits, random perturbations often break syntax or semantics (e.g., undefined variables, etc).

Algorithm. Instead of using random noise for mutation, we use a language model (LM) as a *mutation* operator, leveraging its ability to produce syntactically and semantically valid patches. At each iteration t, the LM generates a batch of patches $\mathcal{Y}^{t+1} = \{y_1^{t+1}, \dots, y_M^{t+1}\}$ conditioned on a set of prior patches \mathcal{E}^t : $\mathcal{Y}^{t+1} \sim \pi(\cdot \mid x, C(x), \mathcal{E}^t)$. We refer to \mathcal{E}^t as conditioning examples consisting of patches generated at iteration t. Following the selection step in evolutionary algorithms, \mathcal{E}^t could be selected as the top-K patches ranked by a scoring function R (i.e., fitness function in evolutionary algorithms). Note that we find that our model after training can self-evolve without this selector (see Section 4.3 and Section 5.2), so this step is optional. The full procedure is detailed in Algorithm 1.

Question: Can a language model naturally perform mutation? Ideally, the mutation operator should generate patches that improve scores. However, as shown in Section 5.2, models trained with classical SFT—conditioned only on the issue and code context—struggle to refine existing patches. In the next section, we present our approach to overcome this limitation.

4.2 SMALL-SCALE MUTATION SUPERVISED FINE-TUNING

Classical supervised fine-tuning (SFT) fails at mutation because it never learns to condition on previous patches. To train the model for mutation, it must observe *conditioning examples*—patches

217

218

219

220

221

222

223

224

225

226

227

228 229

230

231

232 233

235

236

237 238

239

240

241

242

243

244

245

246

247

248

249

250

251 252

253

254

255 256

257 258

259 260 261

262

263

264

265

266

267

268 269

from previous iterations—so it can learn to refine them. In EvoScale, conditioning examples are drawn from the model's earlier outputs. We introduce a two-stage supervised fine-tuning (SFT) process: classical SFT followed by mutation SFT. The classical SFT model is first trained and then used to generate conditioning examples for training the mutation SFT model.

Stage 1 — Classical SFT. We fine-tune a base model on inputs consisting of the issue description xand code context C(x), with targets that include a chain-of-thought (CoT) trace and the ground-truth patch, jointly denoted as y_{SFT}^* . Following prior work on dataset curation (Yang et al., 2025; Xie et al., 2025), we use a teacher model μ (e.g., a larger LLM; see Section 5.1) to generate CoT traces. The training objective is:

$$\max_{\pi_{\mathsf{SFT}}} \mathbb{E}_{x \sim \mathcal{D}, y_{\mathsf{SFT}}^* \sim \mu(\cdot \mid x, C(x))} \left[\log \pi_{\mathsf{SFT}} (y_{\mathsf{SFT}}^* \mid x, C(x)) \right]. \tag{1}$$

We refer to the resulting model π_{SFT} as the classical SFT model.

Stage 2 — Mutation SFT. We fine-tune a second model, initialized from the same base model, using inputs x, C(x), and a set of conditioning examples \mathcal{E} consisting of patches sampled from the classical SFT model $\pi_{\rm SFT}$. The target $y_{
m M-SFT}^*$ includes a CoT trace generated by the teacher model μ conditioned on \mathcal{E} , along with the ground-truth patch. The training objective is:

$$\max_{\pi_{\text{M-SFT}}} \mathbb{E}_{x \sim \mathcal{D}, \mathcal{E} \sim \pi_{\text{SFT}}(\cdot \mid x, C(x)), y_{\text{M-SFT}}^* \sim \mu(\cdot \mid x, C(x), \mathcal{E})} \left[\log \pi_{\text{M-SFT}}(y_{\text{M-SFT}}^* \mid x, C(x), \mathcal{E}) \right]. \tag{2}$$

We refer to the resulting model π_{M-SFT} as the mutation SFT model.

Training on small-scale datasets.

EvoScale targets issues where one-shot generation often fails, but high-scoring patches can still be found through sufficient sampling. This means the model generates a mix of high- and low-scoring patches, so conditioning examples should reflect this diversity. If all examples were already highscoring, test-time scaling would offer limited benefit. Training a classical SFT model on the full dataset, however, leads to memorization, reducing output diversity and making it difficult to construct diverse conditioning examples for mutation SFT. To preserve diversity, we collect y_{SFT}^* and $y_{\text{M-SFT}}^*$ on disjoint subsets of the data. See Appendix D for details.

Limitation of SFT in *self-evolving***.** The mutation SFT model π_{M-SFT} is trained on conditioning examples from the classical SFT model π_{SFT} , which include both low- and high-scoring patches. This raises a natural question: can $\pi_{\text{M-SFT}}$ learn to improve low-scoring patches on its own—i.e., selfevolve—without relying on reward models to select high-scoring examples? If so, we could eliminate the selection step (Line 3 in Algorithm 1), reducing scoring costs and sample usage. However, we find that SFT alone cannot enable self-evolution. Section 4.3 introduces a reinforcement learning approach that trains the model to self-evolve without scoring or filtering.

Algorithm 1 Evolutionary Test-Time Scaling (EvoScale)

Require: Issue description x, code context C(x), editor model π , number of iterations T, samples per iteration M, optional selection size K

- 1: Generate initial outputs $\mathcal{Y}^0 := \{y_1^0, \cdots, y_M^0\} \sim \pi(\cdot \mid x, C(x))$
- 2: **for** t = 1 to T **do**
- (Optional) Select conditioning examples $\mathcal{E}^{t-1} := \{\bar{y}_1^{t-1}, \cdots, \bar{y}_K^{t-1}\} = \operatorname{Select}(\mathcal{Y}^{t-1})$ Generate new outputs $\mathcal{Y}^t := \{y_1^t, \cdots, y_M^t\} \sim \pi(\cdot \mid x, C(x), \mathcal{E}^{t-1})\}$
- 4:
- 5: end for

4.3 LEARNING TO SELF-EVOLVE VIA LARGE-SCALE REINFORCEMENT LEARNING (RL)

To self-evolve, the model must generate patches that maximize a scoring function R, given conditioning examples \mathcal{E} from previous patches. This setup naturally aligns with the reinforcement learning (RL) (Sutton and Barto, 2018), where a policy π is optimized to maximize expected rewards over time. Since our goal is to maximize the reward at the final iteration T, a naïve RL objective is:

$$\max_{\pi} \mathbb{E}_{y^t \sim \pi(\cdot \mid x, C(x), \mathcal{E}^{t-1})} \Big[\sum_{t=0}^{T} r_t \Big], \quad \text{where} \quad r_t = \begin{cases} R(x, y^t), & t = T \\ 0, & \text{otherwise} \end{cases}$$
 (3)

This objective focuses solely on maximizing the final reward. However, it presents two key challenges:

(1) rewards are sparse, with feedback only at iteration T, making learning inefficient (Lee et al., 2024; Shen et al., 2025); and (2) generating full T-step trajectories is computationally expensive (Snell et al., 2025).

Potential shaping alleviates sparse rewards. We address the sparse reward challenge using potential-based reward shaping (Ng et al., 1999), where the potential function is defined as $\Phi(y) = R(x, y)$. The potential reward at step t is:

$$r_t = \Phi(y^t) - \Phi(y^{t-1}) = R(x, y^t) - R(x, y^{t-1}). \tag{4}$$

Unlike the naïve formulation (Equation 3), this provides non-zero potential rewards at every step, mitigating the sparse reward challenge. The cumulative potential reward forms a telescoping sum: $\sum_{t=1}^{T} r_t = R(x, y^T) - R(x, y^0).$ Since y^0 is fixed, maximizing this sum is equivalent to maximizing the final score as shown by Ng et al. (1999).

Monotonic improvement via local optimization. While optimizing Equation 3 achieves the optimal final reward, it is computationally expensive due to the need for full T-step trajectories. As a more efficient alternative, we train the model to maximize the potential reward at each individual iteration t (Equation 4), avoiding the cost of generating full T-step trajectories. This local optimization reduces computation and runtime while ensuring monotonic reward improvement (see Section 5.2), which is sufficient for improving patch scores over iterations. We formally show this property in Section 4.4.

Implementation. Using the full dataset, we fine-tune the mutation SFT model $\pi_{\text{M-SFT}}$ to maximize the expected potential rewards (Equation 4) in score between a newly generated patch y and a previous patch y' drawn from the conditioning examples \mathcal{E} :

$$\max_{\pi_{\text{RL}}} \mathbb{E}_{y \sim \pi_{\text{RL}}(\cdot \mid x, C(x), \mathcal{E}), y' \sim \mathcal{E}} \left[R(x, y) - R(x, y') - \lambda F(y) \right]. \tag{5}$$

This objective encourages the model to generate patches that consistently improve upon previous ones. To ensure the outputs follow the required syntax, we incorporate a formatting penalty term F into the reward function (see Appendix D for details). The conditioning patch y' is sampled from conditioning examples constructed using patches generated by earlier models, such as π_{SFT} or intermediate checkpoints of π_{RL} .

4.4 THEORETICAL ANALYSIS

We analyze the RL objective in Equation 5, which leverages potential-based reward shaping (Ng et al., 1999), and show that the induced policy yields non-decreasing scores at each iteration.

Assumption 1 (Φ -monotonicity). Let \mathbb{Y} be the set of all patches and $\Phi \colon \mathbb{Y} \to \mathbb{R}$ a potential function. For every $y \in \mathbb{Y}$, there exists a finite sequence $y = y_0, y_1, \dots, y_k$ such that $\Phi(y_{t+1}) \geq \Phi(y_t)$ for all $0 \leq t < k$.

This ensures that from any initial patch one can reach higher-scoring patches without decreasing Φ . **Definition 1** (Myopic Policy). *Define the one-step action-value* $Q_0(y,y') = \Phi(y') - \Phi(y), \ y,y' \in \mathbb{Y}$. The myopic policy π_0 selects, at each state y, any successor that maximizes Q_0 : $\pi_0(y) \in \arg\max_{y' \in \mathbb{Y}} [\Phi(y') - \Phi(y)]$.

Proposition 1 (Monotonic Improvement). Under Assumption 1, any trajectory $\{y^t\}_{t\geq 0}$ generated by the myopic policy π_0 satisfies $\Phi(y^t) \geq \Phi(y^{t-1})$ and $r_t = \Phi(y^t) - \Phi(y^{t-1}) \geq 0 \quad \forall t \geq 1$.

Proof. By definition of π_0 , at each step $y^t \in \arg\max_{y'} \left[\Phi(y') - \Phi(y^{t-1})\right]$. Hence $\Phi(y^t) - \Phi(y^{t-1}) \geq 0$, which immediately gives $\Phi(y^t) \geq \Phi(y^{t-1})$ and $r_t \geq 0$. In particular, training with the potential reward in Equation equation 5 guarantees that

$$R(x, y^t) = \Phi(y^t) \ge \Phi(y^{t-1}) = R(x, y^{t-1}) \quad \forall t.$$

Thus the learned policy produces non-decreasing scores over iterations.

5 EXPERIMENTS

5.1 SETTINGS

Implementation Details. We adopt a pipeline-based scaffold consisting of a retriever and a code editing model (see Appendix C). Both components are trained using small-scale SFT and large-scale

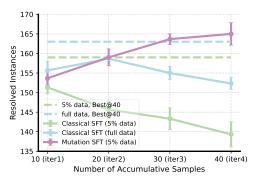
RL. We use the <code>Qwen2.5-Coder-32B-Instruct model Hui et al. (2024)</code> as our base model due to its strong code reasoning capabilities. Our training data is sourced from SWE-Fixer (Xie et al., 2025) and SWE-Gym (Pan et al., 2025). After filtering and deduplication, we obtain a total of 29,404 high-quality instances. For RL training of the code editing model, we rely on a reward model trained on data collected from open source data with 1,889 unique instances. Additional experimental details are provided in Appendix <code>D</code>.

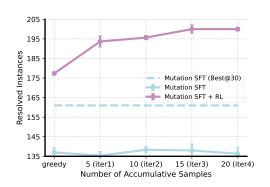
Evaluation and Metrics. We consider two metrics in our evaluation: (1) Greedy: zero-shot pass@1 accuracy, which measures the number of correctly solved instances using greedy generation with syntax retrial (i.e., random sampling up to five times until syntactically correct); (2) Best@N: accuracy of the optimal sample selected by the verifier among N randomly generated samples. Greedy evaluates the model's budget-efficient performance, while Best@N represents the model's potential for test-time scaling.

Test-time Scaling Methods. We evaluate the following test-time scaling methods: (1) Reward Model Selection: selects the optimal patch sample with the highest reward model score; (2) Unit Tests Selection: selects the optimal patch sample based on whether it passes unit tests, including both regression and reproduction tests. If multiple samples pass, one is selected at random; (3) EvoScale: at each evolution iteration, the model generates M patch samples and selects $K \leq M$ samples as the conditional prompt for the next generation. The selection of the K samples is guided by the reward model. In our experiments, we set M=10, K=5, and perform up to four iterations of evolution.

5.2 ANALYSIS

In this section, we present a comprehensive analysis of the proposed EvoScale approach. To simplify our analysis, we use ground-truth localization (retrieval) and focus on the code editing part. All reported results are averaged over three random trials. More results are provided in Appendix A.





(a) RM as selector: Classical SFT v.s. Mutation SFT

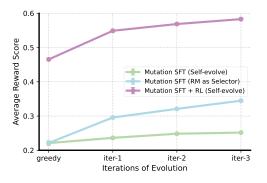
(b) Self-evolve: Mutation SFT v.s. RL

Figure 4: Evolutionary Capability of Different Stages of SFT and RL Models. (a) Reward Model selects the top-5 patch candidates from 10 samples from the previous iteration, and the model iteratively evolves by generating new 10 samples conditioned on the candidates. Performance of the top-1 sample selected by RM is reported. Without the additional mutation SFT training, the model fails to exhibit evolutionary behavior, even when scaling up the training set. (b) Without RM selection, the model only iteratively evolves by conditioning on 5 random samples from the last iteration. RL training improves the model's initial performance and incentivizes the self-evolution capability, while the SFT model fails to self-evolve without guidance from RM.

Can LLMs Iteratively Evolve without Mutation SFT Training? First, we investigate whether the mutation SFT is necessary for LLMs to learn how to iteratively improve their generations. Specifically, we fine-tune base LLMs using either classical SFT (without conditional generation) or mutation SFT. As shown in Figure 4(a), models trained with classical SFT fail to improve their outputs when conditioned on previous samples. In contrast, mutation SFT enables the model to iteratively improve under the guidance of a reward model. The performance of the mutation SFT model at later iterations can surpass the classical SFT model by scaling up the samples (e.g., Best@40). Moreover, this iterative refinement capability can be learned effectively even with a small number of training data.

¹https://huggingface.co/nebius

RL Enables Self-evolve Capability. While mutation SFT model demonstrates evolutionary behavior when guided by a reward model, we further examine whether it can self-evolve without such guidance. Specifically, instead of selecting the top-K candidates to ensure generation quality, we allow the model to generate M=K=5 random samples for the next iteration of conditional generation. However, as shown in Figure 4(b), the SFT model fails to learn self-evolution without reward model selection. Interestingly, RL training significantly improves the SFT model in two key aspects. First, RL substantially boosts the model's greedy performance, surpassing even the Best@N performance of 30 randomly generated samples from the SFT model. Second, we observe that the RL-trained model exhibits strong self-evolution capability: even when conditioned on its random outputs, the model can self-refine and improve performance across iterations without reward model guidance. We provide further analysis of the model's behavior through demo examples in Appendix B.1.



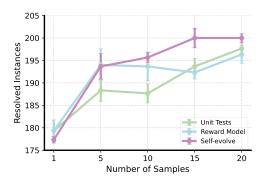


Figure 5: Average Reward Score of Patch Samples at Each Evolution Iteration. Reward scores are normalized via a sigmoid function before average. The SFT model struggles to improve reward scores without the guidance of a reward model to select top-K conditional patch samples, while the RL model consistently self-improves its reward score across iterations without external guidance, validating our theoretical results of monotonic improvement in Section 4.4

Figure 6: Comparison with Other Test-Time Scaling Methods. Reward model selection requires deploying an additional model at test time and can become unstable as the number of samples increases. Unit test selection is computationally expensive and performs poorly with a small sample size. In contrast, self-evolution demonstrates high sample efficiency and strong test-time scaling performance.

Do our SFT and RL Models Monotonically Improve Reward Scores over Iterations? We further analyze the evolutionary behavior of the SFT and RL models by measuring the average reward score of the patch samples generated at each iteration. As shown in Figure 5, although the SFT model learns to iteratively improve reward scores, it relies on the reward model to select high-quality conditioning examples to achieve significant improvements. In contrast, the RL model trained with potential-based reward, naturally learns to self-evolve without any external guidance. Its reward scores improve monotonically across iterations, aligns with our theoretical analysis in Section 4.4.

Evolutionary Test-time Scaling v.s. Other Test-time Scaling Methods. Next, we further compare evolutionary test-time scaling with other test-time scaling methods. Starting from the RL model, we first randomly sample N=5,10,15,20 patch samples and let the reward model and unit tests select the best sample among the subsets. In addition, we let the model perform self-evolution with K=5 samples per iteration, up to four iterations (20 samples in total). The test-time scaling results presented in Figure 6 demonstrate both efficiency and effectiveness of evolutionary test-time scaling.

More Results. We include additional analysis in Appendix A, including different reward modeling methods in RL, the impact of mutation sampling temperature, the impact of RL training, and runtime comparison of different test-time scaling methods.

5.3 RESULTS IN THE WILD: SWE-BENCH PERFORMANCE

We present the main results of our RL-trained model, EvoScale-32B, on the SWE-bench Verified benchmark (Jimenez et al., 2024) and compare its performance against both open-source and proprietary systems. We report results for both greedy decoding and Best@N metrics, using our own retrieval framework (see details of retrieval in Appendix C). For test-time scaling, we apply iterative self-evolution, allowing the RL model to generate M=25 samples per iteration. We observe that the

Table 1: **Results on SWE-bench Verified.** EvoScale-32B outperforms all small-scale models under greedy decoding, while achieving comparable performance with current SOTA SWE-RL with much fewer training data and test-time scaling samples.

Model Scale	Model/Methods	Scaffold	SWE-Verified Resolved Rate
	GPT-4o Yang et al. (2024)	SWE-agent	23.0
	GPT-4o Xia et al. (2024)	Agentless	38.8
	GPT-40 Zhang et al. (2024)	AutoCodeRover	28.8
	GPT-40 Ma et al. (2024)	SWE-SynInfer	31.8
	OpenAI o1 Xia et al. (2024)	Agentless	48.0
Lawas	Claude 3.5 Sonnet Yang et al. (2024)	SWE-agent	33.6
Large	Claude 3.5 Sonnet Wang et al. (2025)	OpenHands	53.0
	Claude 3.5 Sonnet Xia et al. (2024)	Agentless	50.8
	Claude 3.5 Sonnet Zhang et al. (2024)	AutoCodeRover	46.2
	Claude 3.7 Sonnet Anthropic (2025)	SWE-agent	58.2
	DeepSeek-R1 Guo et al. (2025)	Agentless	49.2
	DeepSeek-V3 Liu et al. (2024)	Agentless	42.0
	Lingma-SWE-GPT-7B (Greedy) Ma et al. (2024)	SWE-SynInfer	18.2
	Lingma-SWE-GPT-72B (Greedy) Ma et al. (2024)	SWE-SynInfer	28.8
	SWE-Fixer-72B (Greedy) Xie et al. (2025)	SWE-Fixer	30.2
	SWE-Gym-32B (Greedy) Pan et al. (2025)	OpenHands	20.6
Small	SWE-Gym-32B (Best@16) Pan et al. (2025)	OpenHands	32.0
	Llama3-SWE-RL-70B (Best@80) Wei et al. (2025)	Agentless Mini	37.0
	Llama3-SWE-RL-70B (Best@160) Wei et al. (2025)	Agentless Mini	40.0
	Llama3-SWE-RL-70B (Best@500) Wei et al. (2025)	Agentless Mini	41.0
	EvoScale-32B (Greedy)	EvoScale	35.8
	EvoScale-32B (Best@10)	EvoScale	38.9
	EvoScale-32B (Best@25)	EvoScale	40.2
	EvoScale-32B (Best@50)	EvoScale	41.6

initial iterations produce more diverse candidate patches, while later iterations generate higher-quality, more refined patches. To balance diversity and refinement, we aggregate all generated samples across iterations into a combined pool of N=50 candidates. As discussed in Section 5.2, different verifiers provide complementary strengths. We therefore combine both the reward model and unit tests to select the best patch from the candidate pool.

As shown in Table 1, EvoScale-32B achieves a greedy accuracy of 35.8, outperforming all existing small-scale models under greedy decoding. Additionally, it achieves a Best@50 score of 41.6, matching the performance of the current state-of-the-art Llama3-SWE-RL-70B Wei et al. (2025), which requires Best@500 decoding—incurring over 10× higher sampling cost. It is also worth noting that agent-based methods incur even higher test-time computational cost, as each generation corresponds to a full rollout trajectory with multiple interactions. In contrast, EvoScale-32B achieves state-of-the-art performance with significantly lower inference cost and is trained on fewer than 30K open-source samples, compared to millions of proprietary data used to train Llama3-SWE-RL-70B.

6 Concluding Remarks

We propose Evolutionary Test-time Scaling (EvoScale), a sample-efficient inference-time method that enables small language models to approach the performance of 100B+ parameter models using just 50 code patch samples—without requiring interaction trajectories with the runtime environment. EvoScale opens up a new direction for sample-efficient test-time scaling in real-world software engineering tasks: (1) Evolution improves sample efficiency. Our results show that evolutionary strategies, which iteratively refine generations, can drastically reduce the number of required samples. This contrasts with prior work that primarily focuses on improving verifiers (e.g., reward models, test cases); (2) RL enables self-evolution. We show that reinforcement learning (RL) can train models to refine their outputs without relying on external verifiers at inference. While our current method optimizes local reward differences, future work may explore optimizing cumulative potential rewards over entire trajectories. Compared to Snell et al. (2025), who maintains all prior outputs in the prompt during revision, our method retains only the most recent output—making it more suitable for SWE tasks with long context windows; (3) **Limitations and future work.** This work focuses on a pipeline-based (agentless) setup. Extending EvoScale to agentic settings where models interact with code and runtime environments, remains an interesting future work. While our study focuses on SWE due to its realism and difficulty, we believe the core ideas behind EvoScale could generalize to other agentic tasks, an exciting direction for future work.

REFERENCES

- Anthropic. Introducing claude 3.7 sonnet, 2025., 2025. URL https://www.anthropic.com/claude/sonnet.9
- Antonis Antoniades, Albert Örwall, Kexun Zhang, Yuxi Xie, Anirudh Goyal, and William Yang Wang. SWE-search: Enhancing software agents with monte carlo tree search and iterative refinement. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=G7sIFXugTX. 3
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling, 2024. URL https://arxiv.org/abs/2407.21787.1,2
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021. 1
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 1, 3
- Ryan Ehrlich, Bradley Brown, Jordan Juravsky, Ronald Clark, Christopher Ré, and Azalia Mirhoseini. Codemonkeys: Scaling test-time compute for software engineering. *arXiv preprint arXiv:2501.14723*, 2025. 2
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 9
- Nikolaus Hansen. The cma evolution strategy: A tutorial. arXiv preprint arXiv:1604.00772, 2016. 2
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. 1, 3
- Jian Hu, Xibin Wu, Zilin Zhu, Xianyu, Weixun Wang, Dehao Zhang, and Yu Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*, 2024. 33
- Jian Hu, Jason Klein Liu, and Wei Shen. Reinforce++: An efficient rlhf algorithm with robustness to both prompt and reward models, 2025. URL https://arxiv.org/abs/2501.03262.33, 34, 35
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, Yunlong Feng, Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and Junyang Lin. Qwen2.5-coder technical report, 2024. URL https://arxiv.org/abs/2409.12186.7,33
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL https://openreview.net/forum?id=chfJJYC3iL.1

- Naman Jain, Jaskirat Singh, Manish Shetty, Liang Zheng, Koushik Sen, and Ion Stoica. R2e-gym: Procedural environments and hybrid verifiers for scaling open-weights swe agents. *arXiv preprint arXiv:2504.07164*, 2025b. 2, 3
 - Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNOM66. 1, 2, 8
 - Chi-Chang Lee, Zhang-Wei Hong, and Pulkit Agrawal. Going beyond heuristics by imposing policy improvement as a constraint. *Advances in Neural Information Processing Systems*, 37: 138032–138087, 2024. 6
 - Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023. 1, 3
 - Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint* arXiv:2412.19437, 2024. 9
 - Yingwei Ma, Rongyu Cao, Yongchang Cao, Yue Zhang, Jue Chen, Yibo Liu, Yuchen Liu, Binhua Li, Fei Huang, and Yongbin Li. Lingma swe-gpt: An open development-process-centric language model for automated software improvement. *arXiv* preprint arXiv:2411.00622, 2024. 2, 9
 - Yingwei Ma, Yongbin Li, Yihong Dong, Xue Jiang, Rongyu Cao, Jue Chen, Fei Huang, and Binhua Li. Thinking longer, not larger: Enhancing software engineering agents via scaling test-time compute, 2025. URL https://arxiv.org/abs/2503.23803.2,3
 - Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Icml*, volume 99, pages 278–287, 1999. 2, 6
 - Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with SWE-gym. In *ICLR 2025 Third Workshop on Deep Learning for Code*, 2025. URL https://openreview.net/forum?id=lpFFpTbi9s. 1, 2, 3, 7, 9, 14
 - Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017. 2
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.3
 - Maohao Shen, Soumya Ghosh, Prasanna Sattigeri, Subhro Das, Yuheng Bu, and Gregory Wornell. Reliable gradient-free and likelihood-free prompt tuning. In *Findings of the Association for Computational Linguistics: EACL 2023*. Association for Computational Linguistics, 2023. URL https://aclanthology.org/2023.findings-eacl.183/. 2
 - Maohao Shen, Guangtao Zeng, Zhenting Qi, Zhang-Wei Hong, Zhenfang Chen, Wei Lu, Gregory Wornell, Subhro Das, David Cox, and Chuang Gan. Satori: Reinforcement learning with Chain-of-Action-Thought enhances llm reasoning via autoregressive search. *arXiv preprint arXiv:2502.02508*, 2025. 6
 - Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient RLHF framework. In *Proceedings of the Twentieth European Conference on Computer Systems*. ACM, 2025. 33
 - Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=4FWAwZtd2n. 1, 3, 6, 9

 Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. 2018. 5

Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=OJd3ayDDoF. 9

Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*, 2025. 2, 3, 9, 13, 35

Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014. 2

Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents, 2024. URL https://arxiv.org/abs/2407.01489. 1, 2, 3, 9, 14, 31, 32, 35

Chengxing Xie, Bowen Li, Chang Gao, He Du, Wai Lam, Difan Zou, and Kai Chen. Swe-fixer: Training open-source llms for effective and efficient github issue resolution, 2025. URL https://arxiv.org/abs/2501.05040. 1, 2, 3, 5, 7, 9

John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=mXpq6ut8J3. 2, 3, 9

John Yang, Kilian Leret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents, 2025. URL https://arxiv.org/abs/2504.21798. 1, 2, 5

Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 1592–1604, 2024. 9

APPENDIX

A	Add	itional Experiments				
В	Demo Examples					
	B.1	Type 1: Prior patches are all wrong				
	B.2	Type 2: Prior patches are partially wrong				
	B.3					
C	Scaf	fold of EvoScale				
	C .1	Retriever				
	C.2	Code Editing Model				
	C.3	Verifier				
D	Imp	lementation Details				
	$D.\bar{1}$	Dataset Collection				
	D.2	Training Pipeline and Hardware				
	D.3	Retrieval Model				
	D.4	Retrieval Reward Model				
	D.5	Code Editing Model				
	D.6	Code Editing Reward Model				
	D.7					
E	Pro	nnt Template				

DISCLAIMER ON THE USE OF LARGE LANGUAGE MODELS (LLMS)

In preparing this manuscript, Large Language Models (LLMs) were used exclusively for polishing the writing. They were not involved in research ideation, methodology, experiments, analysis, or paper content generation. The authors take full responsibility for all claims and results presented in this paper.

A ADDITIONAL EXPERIMENTS

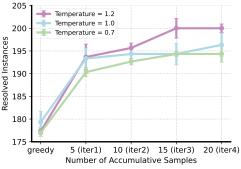
In this section, we provide additional analytical experiments of EvoScale and model training.

RL Reward Modeling: Reward Model is More Effective than String-matching. A reliable reward signal is the key to drive RL training. To better understand the impact of different components in reward modeling, we conduct an ablation study comparing three variants: using only the reward model score, using only string-matching rewards proposed in (Wei et al., 2025), and using both. As shown in Table 2, models trained with a single reward component show degraded greedy decoding performance compared to the model trained with the hybrid reward. In particular, the reward model plays a crucial role in boosting the performance, while the string-matching reward helps the model learn better syntactic structure. However, the results also suggest that naïve string-matching Wei et al. (2025) alone may not serve as a reliable reward signal for SWE tasks.

Table 2: **Ablation Study on Reward Modeling.** The total number of instances is 500. Compared to the SFT model, RL using the RM reward significantly improves performance but introduces more syntax errors. In contrast, RL with a string-matching reward reduces syntax errors but fails to improve reasoning capability. A hybrid reward signal effectively balances both aspects, achieving superior performance.

Metrics	Mutation SFT (5% data)	RM RL	String-matching RL	Hybrid Reward RL
Num of Resolved Instances	137 ± 2.5	171 ± 1.7	140.7 ± 0.5	179.3 ± 2.4
Num of Syntax-correct Instances	427	404	478	471

EvoScale Prefers Higher Mutation Sampling Temperature. Mutation sampling plays a critical role in Evolutionary Test-Time Scaling. To investigate its impact, we vary the model's sampling temperature across 0.7, 1.0, 1.2 and perform self-evolution over four iterations. As shown in Figure 7, higher temperatures demostrate better performance. Intuitively, a larger temperature increases the diversity of generated patch samples, providing richer information for the mutation operator to produce improved patches in subsequent iterations. In contrast, lower temperatures tend to result in repetitive patch samples and may lead the model to converge quickly to suboptimal solutions.



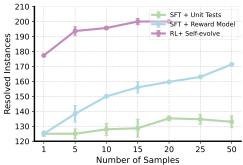


Figure 7: **Impact of Mutation Sampling Temperature.** Higher sampling temperatures in EvoScale encourage greater diversity among mutation samples, leading to more effective iterative improvements.

Figure 8: Classical SFT + Test-time Scaling v.s. Mutation RL + Self-evolve. RL Model with self-evolve capability is more effective than classical SFT model using other test-time scaling methods.

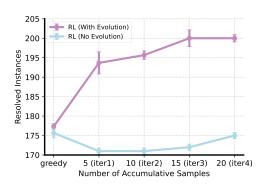
SFT + Test-time Scaling v.s. RL + Self-evolve. In Figure 6, we demonstrated the superior performance of Evolutionary Test-Time Scaling using our RL-trained model. To further investigate this result, we compare against other test-time scaling methods applied to a classical SFT model trained on the full dataset (30K instances), since the typical procedure in most existing SWE work Pan et al. (2025); Xia et al. (2024) trains a SFT model and applying verifiers (e.g., reward models or unit tests) for test-time scaling.

However, as shown in Figure 8, this approach proves to be less effective: (1) With 50 samples, the SFT model's Best@50 performance is still outperformed by the greedy decoding of the RL model, despite both being trained on the same dataset. (2) The SFT model is relatively sensitive to the choice of verifier. When using unit tests (including both reproduction and regression tests) as the verifier, increasing the number of samples results in only marginal performance gains. These observations support our hypothesis: while correct solutions do exist in the SFT model's output distribution, they are rarely sampled due to its dispersed sampling distribution. In contrast, the RL model learns to self-refine the sampling distribution towards high-scoring region.

Table 3: Average Runtime per Instance for Different Test-Time Scaling Methods. Runtime is measured using a sample budget of 10. EvoScale achieves the highest efficiency, while unit test-based selection incurs over 6× higher runtime cost.

Metrics	Unit Tests	Reward Model	Self-evolve
Wall-clock Time (seconds)	92.8 ± 2.6	18.1 ± 0.3	16.6 ± 0.4

Runtime Comparison of Different Test-time Scaling Methods. To evaluate the efficiency of different test-time scaling methods, we measure the average runtime per instance using a sample budget of 10. For our proposed EvoScale approach, the runtime consists solely of iteratively prompting the RL model to generate 10 samples. Reward model selection incurs additional computational cost due to running the reward model to score each sample, and unit test selection requires executing each patch in a sandbox environment. Although unit test selection is effective when scaling to larger sample sizes (see Figure 6), it comes at a cost around 6× slower than EvoScale.



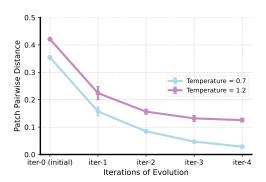


Figure 9: **RL with vs. without Self-Evolution Training.** Removing evolution training during the RL stage results in a model that lacks iterative self-improvement capabilities.

Figure 10: **Diversity of Generated Patches over Iterations.** The generated patches gradually converge toward a high-scoring distribution, while higher temperatures encourage greater diversity.

Would RL without Evolution Training still Work? We consider a simplified training setup for the code editing model, where the base model is trained using classical SFT followed by RL without incorporating mutation data or potential-based rewards. As shown in Figure 2, although this simplified RL approach can still improve the SFT model's greedy performance, it fails to equip the model with iterative self-improvement ability. This finding demonstrates the importance of evolution training, particularly the use of potential-based rewards, in incentivizing the model to learn how to self-refine over multiple iterations.

Analyzing Diversity of Generated Patches. To better understand the behavior of the evolutionary generation, we measure how the generated patch diversity changes over iterations. To quantify this, we compute the average pairwise distance between the generated patches for each instance and report the mean over the entire dataset. Lower values indicate reduced diversity, i.e., more similar patches. We evaluate our RL model at two sampling temperatures (0.7 and 1.2). As shown in Figure 10, diversity decreases over iterations, which aligns with our goal of gradually converging towards high-quality solutions (illustrated in Fig. 3). However, rapid convergence risks local optima. As shown in Figure 7, using higher temperatures helps preserve diversity and can increase performance.

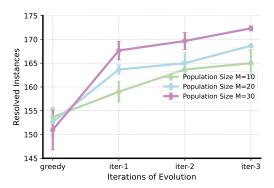


Figure 11: Performance of Mutation SFT (RM as Selector) vs. Population Size (M). Increasing the population size raises the performance upper bound through evolution, but also results in higher inference costs.

Impact of Population Size. While our RL model can self-evolve without relying on a reward model, it is important to note that population size plays an important role only when a selection mechanism, such as a reward model, is available to choose among candidates. To investigate this, we conduct the experiment shown in Figure 4(a) using the Mutation SFT model. Figure 11 reports the Best-of-N (BoN) performance of the SFT model under different population sizes (M=10,20,30). We observe that larger population sizes yield higher performance upper bounds after multiple evolution iterations,

but at the price of increased inference costs. Therefore, we adopt M=10 in our main experiments, as it achieves strong performance after four iterations while maintaining relatively low cost.

B DEMO EXAMPLES

The core idea of EvoScale is to use a mix of correct and incorrect patches as context to generate a correct patch. To do this effectively, the model must demonstrate the following capabilities:

- Generate correct patches even when all prior patches are incorrect—essential for making progress when initial generations fail.
- 2. Generate correct patches when most prior patches are incorrect—common in practice, requiring the model to identify and build on the few correct ones.
- 3. Preserve correctness when all prior patches are correct—ensuring the model doesn't over-edit and degrade already correct solutions.

We present examples below to illustrate each of these capabilities.

B.1 Type 1: Prior patches are all wrong

```
problem Statement 1

sum: min_count is not available for reduction with more than one dimensions
**Is your feature request related to a problem? Please describe.**

'sum' with 'min_count' errors when passing more than one dim:

'``python
import xarray as xr
da = xr.DataArray([[1., 2, 3], [4, 5, 6]])
da.sum(["dim_0", "dim_1"], min_count=1)

**Describe the solution you'd like**
The logic to calculate the number of valid elements is here:
https://github.com/pydata/xarray/blob/lbe777fe725a85b8cc0f65a2bc4lf4bc2bal8043/xarray/
core/nanops.py#L35

I *think* this can be fixed by replacing

'mask.shape[axis]' with 'np.take(a.shape, axis).prod()'

**Additional context**
nongPotentially relevant for #4351
```

```
Conditional patch 1 (Wrong)
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..3e6c5db 100644
--- a/xarray/core/nanops.pv
+++ b/xarray/core/nanops.py
@@ -32,7 +32,16 @@ def _maybe_null_out(result, axis, mask, min_count=1):
     if axis is not None and getattr(result, "ndim", False):
         null_mask = (mask.shape[axis] - mask.sum(axis) - min_count) < 0</pre>
         if np.isscalar(axis):
              # in the scalar case we can just use the shape of that axis
             size = mask.shape[axis]
         else:
              \ensuremath{\text{\#}} in the case of multiple axis we need to take the product
              \ensuremath{\text{\#}} of all the axis sizes since we want the total number of
              \ensuremath{\sharp} elements we are reducing, not just the size of one axis
              size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
              dtype, fill_value = dtypes.maybe_promote(result.dtype)
              result = result.astype(dtype)
```

865 866

867

868

869 870

871

872

873

874

875

876 877

878

879

880 881 882

883 884

885

886

887 888

889

890

891

892

893 894

895

896

897

898 899 900

901 902

903

904

905 906

907

908

909

910

911

912

913

914

915

```
Conditional patch 2 (Wrong)
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..f533eed 100644
   a/xarray/core/nanops.py
+++ b/xarray/core/nanops.py
@@ -32,7 +32,15 @@ def _maybe_null_out(result, axis, mask, min_count=1):
    if axis is not None and getattr(result, "ndim", False):
         null_{mask} = (mask.shape[axis] - mask.sum(axis) - min_count) < 0
         # handle the case where we are reducing on multiple axis
         # we need to consider the product of their sizes, not just any
         # one of them.
        if np.isscalar(axis):
             size = mask.shape[axis]
             size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
             dtype, fill_value = dtypes.maybe_promote(result.dtype)
             result = result.astype(dtype)
```

```
Conditional patch 3 (Wrong)
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..f533eed 100644
--- a/xarray/core/nanops.py
+++ b/xarray/core/nanops.py
@@ -32,7 +32,15 @@ def _maybe_null_out(result, axis, mask, min_count=1):
     if axis is not None and getattr(result, "ndim", False):
         null_mask = (mask.shape[axis] - mask.sum(axis) - min_count) < 0
# handle the case where we are reducing on multiple axis</pre>
         # we need to consider the product of their sizes, not just any
         # one of them.
         if np.isscalar(axis):
              size = mask.shape[axis]
              size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
              dtype, fill_value = dtypes.maybe_promote(result.dtype)
              result = result.astype(dtype)
```

```
Conditional patch 4 (Wrong)
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..f533eed 100644
--- a/xarray/core/nanops.py
+++ b/xarray/core/nanops.py
@@ -32,7 +32,15 @@ def _maybe_null_out(result, axis, mask, min_count=1):
     if axis is not None and getattr(result, "ndim", False):
         null_mask = (mask.shape[axis] - mask.sum(axis) - min_count) < 0
# handle the case where we are reducing on multiple axis</pre>
          # we need to consider the product of their sizes, not just any
          # one of them.
         if np.isscalar(axis):
              size = mask.shape[axis]
         else:
              size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
              dtype, fill_value = dtypes.maybe_promote(result.dtype)
              result = result.astype(dtype)
```

919

921

922

923 924

925

926

927

928

929 930

931

932

933

934 935

936 937

938

939

940

941

942

943

944

945

946

947

948

949

951

952

953

954

955 956

957

958

959 960

961 962

963

964 965

966

967

968

969 970

971

```
Conditional patch 5 (Wrong)
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..f533eed 100644
   a/xarray/core/nanops.py
+++ b/xarray/core/nanops.py
@@ -32,7 +32,15 @@ def _maybe_null_out(result, axis, mask, min_count=1):
     if axis is not None and getattr(result, "ndim", False):
         null_{mask} = (mask.shape[axis] - mask.sum(axis) - min_count) < 0
         # handle the case where we are reducing on multiple axis
         # we need to consider the product of their sizes, not just any
         # one of them.
        if np.isscalar(axis):
             size = mask.shape[axis]
             size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
             dtype, fill_value = dtypes.maybe_promote(result.dtype)
             result = result.astype(dtype)
```

Thinking Process 1. Peer Patch Review: - All five patches are essentially identical in their approach to solving the issue. They \hookrightarrow all modify the same section of code in _maybe_null_out() to handle the case where \hookrightarrow axis is a list/tuple by calculating the product of the sizes of the reduced dimensions instead of just using the size of a single axis. - The key insight is correct - when reducing over multiple dimensions, the total number \hookrightarrow of elements being reduced is the product of the sizes of those dimensions, not just \hookrightarrow the size of one dimension. - All patches correctly identify that when axis is a scalar, we should use the single \hookrightarrow dimension size as before, but when axis is a list/tuple, we should use the product of the sizes of the specified dimensions. - The solution is consistent with the issue description which suggests using \hookrightarrow np.take(a.shape, axis).prod() to get the total number of elements being reduced. 2. Issue Analysis: - The problem occurs when using sum() with min_count on multiple dimensions because the \hookrightarrow current code only checks the size of one dimension (the first in the axis list) → rather than the total number of elements being reduced. - The intended behavior is that min_count should consider the total number of valid \hookrightarrow elements across all reduced dimensions, not just the number of elements along a \hookrightarrow single dimension. - The fix requires modifying the logic that calculates the number of valid elements when → reducing over multiple dimensions to use the product of the sizes of all reduced \hookrightarrow dimensions. 3. Solution Approach: - The key change needed is in the maybe null out() function where it calculates → null mask for multi-dimensional reductions. - When axis is a list/tuple, we should calculate the total number of elements being \hookrightarrow reduced by taking the product of the sizes of the specified dimensions (using → np.take(mask.shape, axis).prod()). - For single dimension reductions, we can keep the existing behavior of using just → mask.shape[axis]. - The rest of the function's logic (checking min_count against the number of valid \hookrightarrow elements) can remain the same. 4. Code Localization:

- The issue is in xarray/core/nanops.py in the _maybe_null_out() function.

```
- The specific location is where it calculates null_mask for multi-dimensional reductions

→ (when axis is not None and result has ndim).

- The current code uses mask.shape[axis] which is incorrect for multi-dimensional

→ reductions. We need to modify this to use the product of the sizes of all reduced

→ dimensions when axis is a list/tuple.
```

```
Evolved Patch
diff --git a/xarray/core/nanops.py b/xarray/core/nanops.py
index 74e25c7..c642ed6 100644
  - a/xarray/core/nanops.py
+++ b/xarray/core/nanops.py
@@ -27,12 +27,22 @@ def _maybe_null_out(result, axis, mask, min_count=1):
     xarray version of pandas.core.nanops._maybe_null_out
     if hasattr(axis, "__len__"): # if tuple or list
         raise ValueError(
              "min_count is not available for reduction with more than one dimensions."
         # we still need to handle the case where axis is a list/tuple
         # by calculating the product of the sizes of the reduced
         \ensuremath{\text{\#}} dimensions rather than just using one dimension's size
         pass
     if axis is not None and getattr(result, "ndim", False):
         null_mask = (mask.shape[axis] - mask.sum(axis) - min_count) < 0</pre>
         if np.isscalar(axis):
             # in the scalar case we can just use the shape of that axis
             size = mask.shape[axis]
         else:
              \ensuremath{\sharp} in the case of multiple axis we need to take the product
              # of all the axis sizes since we want the total number of
              \ensuremath{\sharp} elements we are reducing, not just the size of one axis
             size = np.take(mask.shape, axis).prod()
         null_mask = (size - mask.sum(axis) - min_count) < 0</pre>
         if null_mask.any():
             dtype, fill_value = dtypes.maybe_promote(result.dtype)
              result = result.astype(dtype)
```

B.2 Type 2: Prior patches are partially wrong

```
Problem Statement 1

Incorrect parenthesizing of Subs
Here is an example.
'''python
>>> from sympy import Subs
>>> from sympy.abc import x,y
>>> 3*Subs(-x+y, (x,),(1,))

LaTeX printing of this gives:
'''python
'3 \\left. - x + y \\right|_{{\substack{ x=1 }}'

![image] (https://quicklatex.com/cache3/76/ql_9672fd7e62c909ff3d9ac8543c2e2576_13.png)

It would be better to be parenthesized to:
'''python
'3 \\left. \\left(- x + y\\right) \\right|_{{\substack{ x=1 }}'

![image] (https://quicklatex.com/cache3/bf/ql_936ffdb876e784206d4c54bb93d28dbf_13.png)
```

```
Conditional patch 1 (Wrong)

diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
index 675a3c4..8b39416 100644
--- a/sympy/printing/latex.py
```

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050 1051 1052

1053

1054

1055

1056

1057

1058

1059

1061

1062

1063

```
1026
           +++ b/sympy/printing/latex.py
1027
          @@ -698,7 +698,12 @@ class LatexPrinter(Printer):
1028
               def _print_Subs(self, subs):
1029
                    expr, old, new = subs.args
1030
                    latex_expr = self._print(expr)
                    # get the expression in proper parenthesized form
1031
                       note: we need to use a lower precedence than MUL
1032
                        to get parens for things like 3*Subs(3+2,...) \longrightarrow 3*(3+2)|...
                        the MUL prec level is the highest that always adds the \,
1033
                        needed parens
1034
                    latex_expr = self.parenthesize(expr, PRECEDENCE["Mul"] - 1)
                    latex_old = (self._print(e) for e in old)
1035
                    latex_new = (self._print(e) for e in new)
1036
                    latex_subs = r' \ '.join(
1037
1038
```

```
Conditional patch 2 (Correct)

diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
index 675a3c4..2d262c4 100644
--- a/sympy/printing/latex.py
+++ b/sympy/printing/latex.py
@@ -698,7 +698,10 @@ class LatexPrinter(Printer):

def _print_Subs(self, subs):
    expr, old, new = subs.args
- latex_expr = self._print(expr)
+    # issue 12165: need to parenthesize the expression if it is
+    # an Add (or lower precedence) to get proper grouping
+    # when a coefficient multiplies a substituted expression
+ latex_expr = self.parenthesize(expr, PRECEDENCE["Mul"])
    latex_old = (self._print(e) for e in old)
    latex_new = (self._print(e) for e in new)
    latex_subs = r'\\ '.join(
```

```
Conditional patch 3 (Correct)
diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
index 675a3c4..b5b4285 100644
  a/sympy/printing/latex.py
+++ b/sympy/printing/latex.pg
@@ -698,7 +698,11 @@ class LatexPrinter(Printer):
     def _print_Subs(self, subs):
         expr, old, new = subs.args
         latex_expr = self._print(expr)
         # always parenthesize the expr so Subs is properly grouped
         \# when multiplied, e.g. 3*Subs(-x+y, (x,), (1,)) vs 3*(-x+y)|_x=1
         # we use Mul precedence since we want to group it at the same
         # level as args to Mul (but not as high as Pow which wouldn't group)
         latex_expr = self.parenthesize(expr, PRECEDENCE["Mul"])
         latex_old = (self._print(e) for e in old)
         latex_new = (self._print(e) for e in new)
latex_subs = r'\\ '.join(
```

```
1066
          Conditional patch 4 (Wrong)
1067
1068
          diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
          index 675a3c4..868ab44 100644
1069

    a/sympy/printing/latex.py

1070
          +++ b/sympy/printing/latex.py
1071
          @@ -698,7 +698,13 @@ class LatexPrinter(Printer):
1072
               def _print_Subs(self, subs):
1073
                   expr, old, new = subs.args
                    latex_expr = self._print(expr)
1074
                    # we want to always parenthesize the expression, otherwise
1075
                    \# something like 3*Subs(-x+y, (x,), (1,)) will have
                    # incorrect LaTeX output
1076
                    \# to see the issue, compare 3*Subs(-x+y, (x,),(1,)) with
                    \# 3*(-x+y).subs(x, 1) in IPython with init_printing()
1077
                    # so we sub in a low precedence level to be safe
1078
                    latex_expr = self.parenthesize(expr, 0)
                   latex_old = (self._print(e) for e in old)
1079
```

```
1080
                    latex_new = (self._print(e) for e in new)
latex_subs = r'\\ '.join(
1081
1082
1083
           Conditional patch 5 (Wrong)
1084
1085
           diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
           index 675a3c4..dab5b75 100644
1086
           --- a/sympy/printing/latex.py
1087
           +++ b/sympy/printing/latex.py
           @@ -698,7 +698,13 @@ class LatexPrinter(Printer):
1088
1089
                def _print_Subs(self, subs):
                     expr, old, new = subs.args
1090
                     latex_expr = self._print(expr)
1091
                     \# issue 15210: always parenthesize the expression because
1092
                     # otherwise the wrong precedence may give incorrect results
1093
                     \# For example 3*Subs(x+y, x, 2) would be 3*x+y|_{x=2} rather
                     # than the correct 3*(x+y)|_{x=2}
1094
                     latex_expr = self.parenthesize(expr, 0)
1095
                     latex_old = (self._print(e) for e in old)
1096
                    latex_new = (self._print(e) for e in new)
latex_subs = r'\\ '.join(
1097
1098
1099
          Thinking Process
1100
           1. Peer Patch Review:
1101
1102
           - Patch 1:
1103
             - Correctly identifies the need to use a lower precedence than MUL to force parentheses
             - The solution is reasonable since MUL precedence is high enough to always add needed
1104
             - However, using PRECEDENCE["Mul"] - 1 might be too low since we want to group at the
1105
             \hookrightarrow same level as MUL args
1106
           - Patch 2:
1107
             - Uses PRECEDENCE["Mul"] which is a reasonable choice since we want to group at the
1108
             \hookrightarrow same level as MUL args
              - Correctly identifies the need to parenthesize Add expressions when multiplied
1109
             - Simple and correct solution that addresses the core issue
1110
1111
             - Similar to patch 2 but with clearer explanation of why Mul precedence is appropriate
1112
             - Also correctly identifies the need to always parenthesize the expression
             - Good solution that matches the intended behavior
1113
1114
             - Uses precedence 0 which is too low - this would always add parentheses even when not
1115
             \hookrightarrow needed
1116
             - While it would fix the issue, it would also introduce unnecessary parentheses in many
             1117
             - Not the optimal solution since it's too aggressive with parenthesizing
1118
           - Patch 5:
1119
             - Similar to patch 4 but with a clearer explanation of the issue
1120
             - Also uses precedence 0 which has the same problem as patch 4
             - While the intent is correct, the solution is too broad in its application
1121
1122
           2. Issue Analysis:
           - The issue shows that when a Subs expression is multiplied by a coefficient (like \hookrightarrow 3*Subs(-x+y, (x,),(1,))), the expression inside Subs isn't properly parenthesized in
1123
1124

→ LaTeX output

           - The current code just prints the expression directly without considering its precedence
1125
           \hookrightarrow relative to multiplication
1126
           - The correct behavior should always parenthesize the Subs expression when it appears in

→ a context where precedence matters (like multiplication)

1127
1128
           3. Solution Approach:
           - The key is to modify the \_print\_Subs method to always parenthesize the expression using
1129
           \hookrightarrow the MUL precedence level
1130
           - This ensures that when Subs appears in a multiplication context, it will be properly
           \hookrightarrow grouped with parentheses
1131
           - Using MUL precedence is appropriate because we want it to group at the same level as
1132
           \hookrightarrow other MUL arguments but not higher than necessary
1133
```

```
1134
          4. Code Localization:
1135
            The issue is in sympy/printing/latex.py in the _print_Subs method
1136
          - The current implementation directly prints the expression without any parenthesizing
          - We need to modify it to use the parenthesize method with MUL precedence
1137
1138
1139
          Evolved Patch
1140
1141
          diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
          index 675a3c4..6dd89a9 100644
1142
              a/sympy/printing/latex.py
1143
           +++ b/sympy/printing/latex.py
          @@ -698,7 +698,14 @@ class LatexPrinter(Printer):
1144
1145
               def _print_Subs(self, subs):
                    expr. old. new = subs.args
1146
                   latex_expr = self._print(expr)
                    # issue 15210: always parenthesize the expression because
1147
                    # otherwise the wrong precedence may give incorrect results
1148
                    # For example 3*Subs(x+y, x, 2) would be 3*x+y \mid \{x=2\} rather
1149
                    \# than the correct 3*(x+y)|_{x=2}
                    # Using MUL precedence ensures we group at the same level as
1150
                   # other MUL args but not higher than necessary (we don't want
1151
                    # to group for Pow since Pow has higher precedence)
                    latex_expr = self.parenthesize(expr, PRECEDENCE["Mul"])
1152
                   latex_old = (self._print(e) for e in old)
                   latex_new = (self._print(e) for e in new)
latex_subs = r'\\ '.join(
1153
1154
1155
          Problem Statement 2
1156
1157
          Contains.as_set returns Contains
1158
          >>> Contains(x, Reals).as_set()
1159
          Contains(x, Reals)
1160
1161
          This is wrong because Contains is not a set (it's a boolean). It results in failures in
1162

→ other places because it doesn't have as_relational (since it isn't a set). For

          → instance, from https://github.com/sympy/sympy/pull/14965#discussion_r205281989
1163
1164
          >>> Piecewise((6, Contains(x, Reals)), (7, True))
1165
          Traceback (most recent call last):
            File "<stdin>", line 1, in <module>
1166
            File "./sympy/functions/elementary/piecewise.py", line 136, in __new__
1167
              r = cls.eval(*newargs)
1168
            File "./sympy/functions/elementary/piecewise.py", line 185, in eval
              c = c.as set().as relational(x)
1169
          AttributeError: 'Contains' object has no attribute 'as_relational'
1170
1171
1172
          Conditional patch 1 (Correct)
1173
          diff --git a/sympy/sets/contains.py b/sympy/sets/contains.py
1174
          index 0b1e933..3fd780b 100644
1175
           --- a/sympy/sets/contains.py
          +++ b/sympy/sets/contains.py
1176
          @@ -45,4 +45,6 @@ class Contains(BooleanFunction):
1177
               def as_set(self):
                    raise NotImplementedError()
1178
                    \ensuremath{\sharp} The set being contained in is the only set we can return
1179
                    # from a Contains object in any meaningful way
                    return self.args[1]
1180
1181
          Conditional patch 2 (Correct)
1182
1183
          diff --git a/sympy/sets/contains.py b/sympy/sets/contains.py
1184
          index 0b1e933..4c089aa 100644
              a/sympy/sets/contains.py
1185
          +++ b/sympy/sets/contains.py
1186
          @@ -45,4 +45,5 @@ class Contains(BooleanFunction):
```

def as_set(self):

```
1188
                    raise NotImplementedError()
1189
                     # The set that is being tested for contains
1190
                    return self.args[1]
1191
1192
           Conditional patch 3 (Correct)
1193
           diff --git a/sympy/sets/contains.py b/sympy/sets/contains.py
1194
           index 0b1e933..985c88e 100644
1195
           --- a/sympy/sets/contains.py
           +++ b/sympy/sets/contains.py
1196
           @@ -45,4 +45,6 @@ class Contains(BooleanFunction):
1197
                def as set(self):
                    raise NotImplementedError()
1198
                     # The set being contained in is the set representation
1199
                    # of Contains objects
                    return self.args[1]
1200
1201
           Conditional patch 4 (Correct)
1202
1203
           diff --git a/sympy/sets/contains.py b/sympy/sets/contains.py
1204
           index 0b1e933..e732791 100644
           --- a/sympy/sets/contains.py
1205
           +++ b/sympy/sets/contains.py
1206
           @@ -45,4 +45,6 @@ class Contains(BooleanFunction):
                def as set(self):
1207
                    raise NotImplementedError()
1208
                    # Contains(x, S) means "x is in S", so as_set returns S
                    return self.args[1]
1209
1210
1211
           Conditional patch 5 (Wrong)
1212
           diff --git a/sympy/sets/contains.py b/sympy/sets/contains.py
1213
           index 0b1e933..4980364 100644
              a/sympy/sets/contains.py
1214
           +++ b/sympy/sets/contains.py
1215
           @@ -45,4 +45,22 @@ class Contains(BooleanFunction):
                def as_set(self):
1216
                    # try to evaluate Contains(x, S)
1217
                    trv:
                        res = self._eval_evalf()
1218
                        if res == S.true:
1219
                            return S.UniversalSet
                        elif res == S.false:
                             return S.EmptySet
1221
                    except (TypeError, AttributeError):
1222
                        pass
                    # fallback
1223
                raise NotImplementedError()
1224
1225
          Thinking Process
1226
           1. Peer Patch Review:
1227
1228
           - Patch 1: This patch correctly identifies that the set being contained in (self.args[1])
           \,\hookrightarrow\, is the most meaningful set representation of a Contains object. The implementation is
1229

→ simple and correct.

1230
           - Patch 2: Similar to Patch 1 but with slightly different wording in the comment. The
1231
           \hookrightarrow functionality is correct but the comment could be more precise about why this is the
1232

→ right approach.

1233
           - Patch 3: This patch adds a note about potential invalid sets, which isn't strictly
1234
           → necessary since the issue is about the basic behavior. The core implementation is
           \,\hookrightarrow\, correct but the comment is overly cautious.
1235
1236
           - Patch 4: This patch provides a clear explanation of why returning the set

→ (self.args[1]) is the correct behavior - because Contains(x, S) means "x is in S", so
1237
           \,\hookrightarrow\, the set representation is just S. This is a correct and well-explained solution.
1238
           - Patch 5: This patch takes a different approach by trying to evaluate the Contains
1239
           \hookrightarrow expression and return either the universal or empty set based on the result. While
1240
```

creative, this doesn't match the simple and correct solution of just returning the set being contained in. The approach is overly complex for this simple issue.

```
1242
1243
           2. Issue Analysis:
              - The issue occurs when Contains.as_set() returns the Contains object itself instead
              \hookrightarrow of a proper set.
1245
              - This causes problems because other code expects a set with as_relational() method.
1246
              - The correct behavior should be to return the set that x is being tested for
              \hookrightarrow containment in (self.args[1]), since that's the only meaningful set representation
1247
              \hookrightarrow of a Contains object.
1248
           3. Solution Approach:
1249
              - The simplest and most correct solution is to modify as_set() to return self.args[1],
1250
              \,\hookrightarrow\, which is the set being tested for containment.
              - This matches the mathematical interpretation that Contains(x, S) means "x is in S",
1251
              \,\hookrightarrow\, so the set representation is just S.
1252
              - No other evaluation or special cases are needed since the set being contained in is
              \hookrightarrow the only meaningful set representation of a Contains object.
1253
1254
           4. Code Localization:
              - The issue is in sympy/sets/contains.py in the Contains class.
1255
              - The problematic method is as_set() which currently raises NotImplementedError.
1256
              - The fix is to modify this method to return self.args[1].
1257
1258
```

B.3 Type 3: Prior patches are all correct

1259 1260

1261

1262

1263

1264

1265

1266 1267 1268

1269 1270

1271

1272

1273 1274

12751276

1277

1278

1279

1280 1281

1282 1283

1284

1285

1286

1287

12881289

1290

1291

1292

1293

1294

```
Problem Statement 1
 _pytest.capture.EncodedFile mode should not include `b` (binary)
 Thanks for submitting an issue!
Here's a quick checklist for what to provide:
 - [x] a detailed description of the bug or suggestion
Exception when youtube-dl logs to pytest captured output. Youtube-dl looks for `b` in
      out.mode `to decide whether to writes `bytes` or `str`. `_pytest.capture.EncodedFile`
 → incorrectly advertises `rb+`, the mode of the underlying stream. Its `write()` method
\hookrightarrow raises an exception when passed `bytes`.
(pytest-issue-ve3) 01:11:48:nlevitt@Internets-Air-2:/tmp$ py.test test.py

→ session starts

platform darwin -- Python 3.7.3, pytest-4.5.0, py-1.8.0, pluggy-0.11.0
 rootdir: /private/tmp
collected 1 item
test.pv F

→ [100%]

→ FATLURES

\hookrightarrow

→ test_foo

\hookrightarrow
    def test foo():
        youtube_dl.YoutubeDL().extract_info('http://example.com/')
```

```
1296
                   test.py:4:
1297
1298
                   \hookrightarrow
1299
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/YoutubeDL.py:796: in extract_info
1300
                         ie_result = ie.extract(url)
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/extractor/common.py:529: in
1301
                   \hookrightarrow extract
1302
                         ie_result = self._real_extract(url)
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/extractor/generic.py:2245: in
1303
                        _real_extract
1304
                          self.to_screen('%s: Requesting header' % video_id)
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/extractor/common.py:913: in
1305
                        to_screen
1306
                          self._downloader.to_screen('[%s] %s' % (self.IE_NAME, msg))
                   {\tt pytest-issue-ve3/lib/python3.7/site-packages/youtube\_dl/YoutubeDL.py:} 502: in to\_screen in the property of the property 
1307
                          return self.to_stdout(message, skip_eol, check_quiet=True)
1308
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/YoutubeDL.py:516: in to_stdout
                          self._write_string(output, self._screen_file)
1309
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/YoutubeDL.py:505: in
1310
                    → _write_string
                          write_string(s, out=out, encoding=self.params.get('encoding'))
1311
                   pytest-issue-ve3/lib/python3.7/site-packages/youtube_dl/utils.py:1496: in write_string
1312
                         out.write(byt)
1313
1314
1315
                   self = <_pytest.capture.EncodedFile object at 0x10df124a8>, obj = b'[generic] example:
1316
                   → Requesting header\n'
1317
                          def write(self, obj):
1318
                                 if isinstance(obj, six.text_type):
                                        obj = obj.encode(self.encoding, "replace")
1319
                                 elif _PY3:
1320
                                       raise TypeError(
                                                "write() argument must be str, not {}".format(type(obj). name )
1321
1322
                                        TypeError: write() argument must be str, not bytes
1323
                   pytest-issue-ve3/lib/python3.7/site-packages/_pytest/capture.py:437: TypeError
1324
                                                                                                                                                        ==== 1 failed in

→ 2.74 seconds

1325
1326
1327
                   - [x] output of `pip list` from the virtual environment you are using
1328
                   Package
                                              Version
1329
1330
                   atomicwrites 1.3.0
                   attrs
                                              19.1.0
1331
                   more-itertools 7.0.0
1332
                   pip
                                              19.1.1
1333
                                              1.8.0
1334
                   pytest
                                             4.5.0
                   setuptools
                                             41.0.1
1335
                   six
                                              1.12.0
1336
                   wcwidth
                                             0.1.7
                   wheel
                                             0.33.4
1337
                   youtube-dl
                                              2019.5.11
1338
1339
                   - [x] pytest and operating system versions
1340
                   This is pytest version 4.5.0, imported from
1341
                   1342
1343
1344
                   macOS 10.14.4 (18E226)
1345
1346
                   - [x] minimal example if possible
1347
1348
                   pip install pytest youtube-dl
                   py.test test.py
1349
```

1358 1359

1360

1361

1362

1363 1364

1365

1366

1367

1368

1369

1370

1371

1372

1373 1374

1375

1376

1377

1379

1381

1382

1383

1384

1386

1387

1389

1390

```
1350
1351
1352
          test.py:
1353
          import youtube_dl
1354
          def test_foo():
              youtube_dl.YoutubeDL().extract_info('http://example.com/')
1355
1356
```

```
Conditional patch 1 (Correct)
diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py
index 6b2d076..aa57c0d 100644
 -- a/src/_pytest/capture.py
+++ b/src/_pytest/capture.py
@@ -448,6 +448,17 @@ class EncodedFile(object):
         return repr(self.buffer)
     def getattr (self, name):
         # when code inspects the mode attribute (like youtube-dl does)
         # we need to hide the 'b' flag from it - we act as a text stream
         # even though we write to a bytes buffer
        if name == "mode":
             mode = getattr(self.buffer, name, None)
             if mode and "b" in mode:
                 # return a copy of the mode without 'b'
                 # this makes e.g. youtube-dl think this is a text stream
                 # and write str instead of bytes
                 return mode.replace("b", "")
         \ensuremath{\sharp} fall back to regular attribute access
         return getattr(object.__getattribute__(self, "buffer"), name)
```

```
Conditional patch 2 (Correct)
          diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py
          index 6b2d076..01e03c8 100644
              a/src/_pytest/capture.py
          +++ b/src/_pytest/capture.py
1378
          @@ -448,6 +448,20 @@ class EncodedFile(object):
                   return repr(self.buffer)
1380
               def __getattr__(self, name):
                   # youtube-dl specifically looks at out.mode to decide whether to
                     write bytes or str. Since we're a text interface over a bytes
                    interface, we should make it look like we're text.
              https://github.com/ytdl-org/youtube-dl/blob/2019.05.11/youtube_dl/utils.py#L1496
                   # See: https://github.com/pytest-dev/pytest/issues/5691
1385
                   if name == "mode":
                       mode = getattr(self.buffer, name, None)
                       if mode and "b" in mode:
                           # remove the 'b' from the mode, but otherwise leave it
                            unchanged. This makes youtube-dl think it can write text
1388
                           # to us when in reality we'll encode it to bytes for the
                           # underlying buffer.
                           mode = "".join(ch for ch in mode if ch != "b")
                       return mode
                   return getattr(object. getattribute (self, "buffer"), name)
1392
```

```
1393
          Conditional patch 3 (Correct)
1394
          diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py
1395
          index 6b2d076..14dfb3c 100644
1396
           --- a/src/_pytest/capture.py
          +++ b/src/_pytest/capture.py
1397
          @@ -447,6 +447,21 @@ class EncodedFile(object):
1398
                   """Ensure that file.name is a string."""
                   return repr(self.buffer)
1400
               @property
               def mode(self):
1401
1402
                   Get the mode of the underlying buffer but remove 'b' from it.
                   This is needed because some code (like youtube-dl) checks
1403
```

1415 1416

1417

1418

1419

1420 1421

1422

1423

1424

1425

14261427

1428

1429

1430

1431 1432

14331434

1435 1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446 1447

1448 1449

```
1404
                       for 'b' in mode to decide whether to write bytes or str.
1405
                       Since this is an encoded (text) file, we need to hide that
1406
                       it's wrapping a binary stream.
1407
                       mode = getattr(self.buffer, "mode", None)
1408
                       if mode and "b" in mode:
    # remove 'b' but keep other flags like '+' or 'r'
    mode = "".join(c for c in mode if c != "b")
1409
1410
                       return mode
1411
                        _getattr__(self, name):
1412
                       return getattr(object.__getattribute__(self, "buffer"), name)
1413
```

Conditional patch 4 (Correct) diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py index 6b2d076..df7fb36 100644 - a/src/_pytest/capture.py +++ b/src/_pytest/capture.py @@ -447,6 +447,22 @@ class EncodedFile(object): """Ensure that file.name is a string.""" return repr(self.buffer) @property def mode(self): """Lie about our mode - remove 'b' if present. Some libraries (like youtube-dl) inspect the mode to decide whether to write bytes or str to the stream. Since we handle that in .write() and .writelines() already, we should pretend to be a text stream even if our underlying buffer is a byte stream. # just in case the buffer doesn't even have a mode... mode = getattr(self.buffer, "mode", "" # ... and remove the 'b' mode = mode.replace("b", "") return mode def __getattr__(self, name): return getattr(object.__getattribute__(self, "buffer"), name)

```
Conditional patch 5 (Correct)
diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py
index 6b2d076..f94e9dc 100644
--- a/src/_pytest/capture.py
+++ b/src/_pytest/capture.py
@@ -447,6 +447,14 @@ class EncodedFile(object):
         """Ensure that file.name is a string."""
         return repr(self.buffer)
     @property
     def mode(self):
         # youtube-dl checks for 'b' in mode to decide whether to write bytes or str.
         # since we are a text interface (we encode strs to bytes), we should not
         # appear as a binary stream, regardless of the underlying buffer's mode.
         # https://github.com/pytest-dev/pytest/issues/5295
         return self.buffer.mode.replace("b", "")
     def __getattr__(self, name):
         return getattr(object.__getattribute__(self, "buffer"), name)
```

```
Thinking Process
1450
1451
          1. Peer Patch Review:
1452
1453
          - Patch 1:
            - Correctly identifies the need to modify the mode attribute when accessed
1454
            - Uses \_\_getattr\_\_ to intercept the mode attribute and remove 'b'
1455
            - Simple and effective solution that matches the issue requirements
            - Correctly maintains all other attribute behavior
1456
1457
          - Patch 2:
```

```
1458
            - Similar to patch 1 but uses a different string manipulation method
1459
            - Also correctly intercepts the mode attribute and removes 'b'
1460
            - Works just as well as patch 1 but with slightly different implementation
1461
1462
            - Takes a different approach by adding a mode property directly
            - This is a cleaner solution since it explicitly defines the behavior
1463
            - Correctly handles the case where the buffer doesn't have a mode
1464
            - Maintains all other attribute behavior through __getattr__
1465
          - Patch 4:
1466
            - Similar to patch 3 but with more detailed docstring
            - Also adds the mode property to remove 'b' from the underlying mode
1467
            - Correctly handles the case where the buffer doesn't have a mode
1468
            - Maintains all other attribute behavior through __getattr__
1469
          - Patch 5:
1470
            - Simplest solution that just adds a mode property
            - Directly returns the underlying mode with 'b' removed
1471
            - Doesn't handle the case where the buffer might not have a mode attribute
1472
            - Still solves the core issue but is less robust than patches 3-4
1473
          2. Issue Analysis:
1474
          - The problem occurs when youtube-dl checks the mode of the captured output stream
          - It sees 'b' in the mode and tries to write bytes when the EncodedFile expects str
1475
          - The solution needs to make the EncodedFile appear as a text stream (no 'b' in mode)
1476
            While still maintaining all other file-like behavior from the underlying buffer
1477
          3. Solution Approach:
1478
          - The key is to intercept access to the 'mode' attribute of the underlying buffer
          - When 'mode' is accessed, we should return the mode string with 'b' removed
1479
          - All other attribute access should continue to work normally from the buffer
1480
          - This matches exactly what youtube-dl expects while maintaining proper functionality
1481
          4. Code Localization:
1482
          - The issue is in the EncodedFile class in src/_pytest/capture.py
          - Need to modify the class to add a mode property that removes 'b' from the underlying
1483
          ⇔ buffer's mode
1484
          - All other behavior should remain unchanged
1485
```

```
1486
          Evolved Patch
1487
1488
          diff --git a/src/_pytest/capture.py b/src/_pytest/capture.py
1489
          index 6b2d076..d3b3473 100644
           --- a/src/_pytest/capture.py
1490
          +++ b/src/_pytest/capture.py
1491
          @@ -447,6 +447,23 @@ class EncodedFile(object):
                    """Ensure that file.name is a string."""
1492
                   return repr(self.buffer)
1493
               @property
1494
               def mode(self):
1495
                     ""When code inspects our mode (like youtube-dl does), lie and
                   remove the 'b' if present. This makes youtube-dl think we're a text
1496
                   stream and write str instead of bytes, which we can then encode to
1497
                   bytes for the underlying buffer.
1498
                   Note: we don't need to handle the case where the buffer doesn't have
1499
                   a mode attribute because in that case, normal attribute access will
                   raise AttributeError like it should.
1500
1501
                   mode = self.buffer.mode
                   if "b" in mode:
1502
                       # remove the 'b' but otherwise preserve the mode
1503
                       mode = mode.replace("b", "")
                   return mode
1504
1505
               def __getattr__(self, name):
                   return getattr(object.__getattribute__(self, "buffer"), name)
1506
1507
```

```
Problem Statement 2

DateFormat.Y() is not zero-padded.
Description
```

1509 1510

```
1512
1513
           The Y specifier for django.utils.dateformat.DateFormat is supposed to always return a
1514
           \hookrightarrow four-digit year padded with zeros. This doesn't seem to be the case for year < 1000.
1515
1516
           Conditional patch 1 (Correct)
1517
           diff --git a/django/utils/dateformat.py b/django/utils/dateformat.py
1518
           index 38cb1c1..823e893 100644
1519
             -- a/django/utils/dateformat.py
           +++ b/django/utils/dateformat.py
1520
           @@ -314,7 +314,9 @@ class DateFormat(TimeFormat):
1521
                def Y(self):
1522
                     "Year, 4 digits; e.g. '1999'"
1523
                     return self.data.year
                     \# We want to zero-pad the year (this is only useful for year < 1000)
1524
                    # so that we get 4 digits.
return '%04d' % self.data.year
1525
1526
                def z(self): """Day of the year, i.e. 1 to 366."""
1527
1528
1529
           Conditional patch 2 (Correct)
1530
           diff --git a/django/utils/dateformat.py b/django/utils/dateformat.py
1531
           index 38cb1c1..15090a4 100644
1532
            -- a/django/utils/dateformat.py
           +++ b/django/utils/dateformat.p
1533
           @@ -314,7 +314,9 @@ class DateFormat(TimeFormat):
1534
1535
                def Y(self):
                     "Year, 4 digits; e.g. '1999'"
1536
                     return self.data.year
                     \# We want to zero-pad the year, so can't just return the number -- \# hence the "%04d" dance.
1537
1538
                    return '%04d' % self.data.year
1539
                def z(self):
1540
                     """Day of the year, i.e. 1 to 366."""
1541
1542
           Conditional patch 3 (Correct)
1543
           diff --qit a/django/utils/dateformat.py b/django/utils/dateformat.py
1544
           index 38cb1c1..abab54f 100644
1545
               a/django/utils/dateformat.py
1546
           +++ b/django/utils/dateformat.pv
           @@ -314,7 +314,8 @@ class DateFormat(TimeFormat):
1547
                def Y(self):
1548
                     "Year, 4 digits; e.g. '1999'"
1549
                     return self.data.year
                     # We want a four-digit year; 0-padding if necessary
1550
                    return '%04d' % self.data.year
1551
                def z(self):
1552
                     """Day of the year, i.e. 1 to 366."""
1553
1554
           Conditional patch 4 (Correct)
1555
1556
           diff --qit a/djanqo/utils/dateformat.py b/djanqo/utils/dateformat.py
1557
           index 38cb1c1..ecf7449 100644
             -- a/django/utils/dateformat.py
1558
           +++ b/django/utils/dateformat.py
1559
           @@ -314,7 +314,11 @@ class DateFormat(TimeFormat):
1560
                def Y(self):
                     "Year, 4 digits; e.g. '1999'"
1561
                     return self.data.year
1562
                     # We hardcode 4 here instead of using a settings value or
                     # self.data.year.__len__() because we want to always return
# a 4-digit year for consistency with PHP's Y behavior.
1563
1564
                     # See also: https://bugs.php.net/bug.php?id=67077
```

return '%04d' % self.data.year

```
1566
1567
               def z(self):
1568
                     ""Day of the year, i.e. 1 to 366."""
1569
1570
          Conditional patch 5 (Correct)
1571
          diff --git a/django/utils/dateformat.py b/django/utils/dateformat.py
1572
          index 38cb1c1..acd0ae0 100644
1573

    a/django/utils/dateformat.py

1574
           +++ b/django/utils/dateformat.py
          @@ -314,7 +314,9 @@ class DateFormat(TimeFormat):
1575
               def Y(self):
1576
                    "Year, 4 digits; e.g. '1999'"
1577
                    return self.data.year
1578
                   # We want to always return a 4-digit year, zero-padded if necessary
                   # This ensures years < 1000 are properly formatted
1579
                   return '%04d' % self.data.year
1580
               def z(self):
1581
                    """Day of the year, i.e. 1 to 366."""
1582
1583
          Thinking Process
1584
1585
          1. Peer Patch Review:
1586
          - Patch 1:
1587
             - Correctly identifies the need to zero-pad the year using '%04d'
             - The comment explains the issue well for years < 1000
1588
            - Simple and correct solution that matches the issue requirements
1589
1590
            - Similar to Patch 1 but with slightly different comment wording
1591
            - Also correctly uses '%04d' to zero-pad the year
            - Correct solution that addresses the issue
1592
1593
          - Patch 3:
             - Uses '%04d' like the other patches but with simpler comment
1594
            - Correct solution that matches the issue requirements
1595
1596
            - Adds more detailed comment explaining the PHP compatibility
1597
            - While the solution is correct, the additional comment about PHP behavior isn't

→ necessary for this Django-specific issue

1598
            - Still a correct solution to the issue
1599
          - Patch 5:
1600
             - Simple and correct solution using '%04d' with clear comment
1601
            - Correctly addresses the issue requirements
1602
          2. Issue Analysis:
1603
           - The issue states that DateFormat.Y() should always return a 4-digit year padded with

→ zeros

1604
          - Currently it returns the raw year number which doesn't zero-pad for years < 1000
1605
          - The fix requires modifying the Y() method to format the year as a 4-digit string with

→ zero-padding

1606
          - This matches the behavior described in the method's docstring which says "4 digits"
1607
          3. Solution Approach:
1608
          - The solution is to modify the Y() method to use string formatting with \ensuremath{\text{``804d'}}
1609
          - This will ensure the year is always returned as a 4-character string with leading zeros
          1610
          - The change is simple and focused on just the Y() method where the issue occurs
1611
          4. Code Localization:
1612
           - The issue is in django/utils/dateformat.py in the DateFormat class
1613
          - The problematic method is Y() which currently just returns self.data.year
          - Need to modify it to return '%04d' % self.data.year instead
1614
1615
1616
          Evolved Patch
1617
          diff --git a/django/utils/dateformat.py b/django/utils/dateformat.py
1618
          index 38cb1c1..d0e66c3 100644
1619
             - a/django/utils/dateformat.py
```

```
+++ b/django/utils/dateformat.py
@@ -314,7 +314,9 @@ class DateFormat(TimeFormat):

def Y(self):
    "Year, 4 digits; e.g. '1999'"
-    return self.data.year
+    # We want to zero-pad the year to ensure we always get
+    # 4 digits, even for years < 1000.
+    return '%04d' % self.data.year

def z(self):
    """Day of the year, i.e. 1 to 366."""</pre>
```

C SCAFFOLD OF EVOSCALE

We adopt a pipeline-based scaffold consisting of two main components: a retriever and a code editing model. The retriever takes either the GitHub repository structure or the contents of individual code files as input and identifies the files most relevant to the given issue. Once retrieval is complete, the full content of the selected code files is passed to the code editing model.

The code editing model receives both the issue description and the retrieved code content as input and generates a patch to resolve the issue. Additionally, there is an optional verifier component, which can be used to select the best patch from a large pool of candidate samples. We describe each component in detail below.

C.1 Retriever

Our retriever is entirely LLM-based and consists of two components: a retrieval model and a retrieval reward model.



Figure 12: **Retrieval Pipeline.** Given the repository's file structure, the retrieval model first selects the top-5 candidate files. These candidates are then re-scored by the retrieval reward model based on file content, and the top-ranked (Top-1) file is returned as the final result.

Retrieval Model The first stage of our retriever uses a retrieval model to identify the top 5 most relevant files based on the repository structure and the GitHub issue description. We adopt the same format as Agentless Xia et al. (2024) to represent the repository's file structure. Given this representation and the issue description, the retrieval model performs a reasoning process and outputs five file paths from the repository. The model is trained using a combination of small-scale supervised fine-tuning (SFT) and large-scale reinforcement learning (RL), see Appendix D.3 for details.

Retrieval Reward Model The retrieval reward model is designed to refine retrieval results in a more fine-grained manner. After the initial top-5 files are retrieved by the retrieval model, the reward model evaluates each one by considering both the file's code content and the issue description. It then outputs a relevance score for each file, and the file with the highest score is selected as the final target for code editing. The retrieval reward model is a classifier-style LLM trained with a binary classification objective, see Appendix D.4 for training details.

C.2 CODE EDITING MODEL

The code editing model receives a prompt formed by concatenating the issue statement with the code content of the retrieved target file. It performs iterative sampling to enable self-evolution during generation.

In the first iteration, given the issue statement and code context, the model generates five diverse responses, each corresponding to a different patch candidate. These five patch candidates are then appended to the input as a conditional prompt for the next iteration of generation. This iterative process allows the model to progressively refine its outputs.

As discussed in Section 4, the code editing model is trained using a combination of small-scale SFT and large-scale RL. Additional training details are provided in Appendix D.5.

C.3 VERIFIER

As discussed in Section 5.2, our code editing model demonstrates the ability to self-evolve by iteratively refining its own generations. While this process improves the quality of patch samples, incorporating external verifiers to select the optimal patch can further boost performance. For software engineering tasks, we consider two primary types of verifiers: an LLM-based reward model and unit tests.

Code Editing Reward Model The code editing reward model is designed to select the best patch from a pool of candidates. It takes as input the retrieved file's code content, the issue description, and a patch candidate in git diff format. The model then outputs a score indicating the quality of the patch. This reward model is implemented as a classifier-based LLM trained with a binary classification objective (see Appendix D.6 for details).

Unit Tests Unit tests consist of two components: (1) Reproduction tests, which validate whether the original GitHub issue can be reproduced and resolved by the patch; (2) Regression tests, which check whether the patch preserves the existing functionality of the codebase. To construct the regression tests, we extract existing test files from the repository using the Agentless Xia et al. (2024) pipeline. For the reproduction tests, we use a trained test generation model that takes the issue description as input and generates tests aimed at reproducing the issue. This reproduction test generator is trained using supervised fine-tuning (see Appendix D.7). For each instance, we sample 100 reproduction tests and retain 5 valid patches to serve as reproduction tests.

Hybrid Verifiers We combine multiple verifiers to select the most promising patch candidates. The selection process is as follows: (1) Regression tests are applied first. Any patch that fails is discarded; (2) Reproduction tests are then executed on the remaining patches. Candidates are ranked based on how many of the five tests they pass. (3) The top-k (k = 2) unique patches are retained per instance. (4) If no patch passes both regression and reproduction tests, we fall back to using all generated candidates without filtering. (5) Finally, among the remaining patches, the code editing reward model is used to select the candidate with the highest reward score for submission.

D IMPLEMENTATION DETAILS

D.1 DATASET COLLECTION

Our primary training data is sourced from SWE-Fixer and SWE-Gym. To ensure data quality, we apply a comprehensive filtering and deduplication process. Specifically, we discard instances that meet any of the following criteria: (1) extremely short or excessively long issue statements; (2) multimodal or non-text content (e.g., images, videos, LaTeX); (3) presence of unrelated external links; (4) inclusion of commit hashes; (5) patches requiring modifications to more than three files. After applying these filters, we obtain 29,404 high-quality training instances, which we use to train both the retrieval and code editing models.

D.2 TRAINING PIPELINE AND HARDWARE

Both the retrieval and code editing models are trained in two stages: (1) small-scale supervised fine-tuning (SFT) for warm-up, and (2) large-scale reinforcement learning (RL) for self-improvement. The SFT trajectories are generated via chain-of-thought prompting using Deepseek-V3-0324. We adopt the Qwen2.5-Coder-32B-Instruct model Hui et al. (2024) as the base model for all components due to its strong code reasoning capabilities. We utilize OpenRLHF Hu et al. (2024) as the training framework for SFT, and use VERL Sheng et al. (2025) as the training framework for RL. All training runs are conducted on NVIDIA H100 GPUs with 80 GB of memory. For evaluation and model inference, we serve models using the sglang framework², employing tensor parallelism with a parallel size of 8 on NVIDIA H100 GPUs.

D.3 RETRIEVAL MODEL

 Dataset Construction To ensure cost-efficient synthetic data generation, we randomly select 1,470 instances (5% of the full dataset) as the small-scale SFT dataset for training the retrieval model.

To train the model to perform step-by-step reasoning and generate relevant file paths conditioned on both the issue statement and the repository structure, we require training data that includes both intermediate reasoning and final retrieved files. However, the raw data only provides the issue descriptions and the ground-truth retrieved files (extracted from the final patch), without any intermediate reasoning.

To address this, we use Deepseek-V3-0324 with a custom retrieval model prompt (see Appendix E) and apply rejection sampling to collect high-quality chain-of-thought (CoT) reasoning traces. Specifically, we check whether the model's top-5 retrieved files include the ground-truth retrieval files. If so, we retain the response as part of our synthetic SFT data.

For each selected instance, we generate one response via greedy decoding and three additional responses using random sampling (temperature = 0.7), as long as they include the correct retrieval files. This results in four responses per instance.

For RL training, we use the remaining 27,598 instances (95% of the dataset), filtering out prompts whose total sequence length exceeds 16,384 tokens. The RL dataset consists of prompts (issue statement + repo structure) as input and the corresponding ground-truth retrieval files as the final answer, without requiring intermediate reasoning.

Supervised Fine-Tuning We perform supervised fine-tuning (SFT) on the Qwen2.5-Coder-32B-Instruct model using the synthetic SFT dataset described above. The prompt template used for training is identical to the one used to construct the synthetic data (see the retrieval model prompt in Appendix E). We train the model using a cosine learning rate scheduler with an initial learning rate of 1e-5. The model is fine-tuned for one epoch with a batch size of 128 and a maximum sequence length of 32,768 tokens.

Reinforcement Learning To further push the limit of the model's retrieval performance, we apply a large-scale reinforcement learning (RL) stage after the SFT stage. After SFT, the model has learned to reason step-by-step and generates a set of candidate retrieval files, denoted as $\mathcal{Y} = \{y_1, y_2, \dots, y_k\}$, where k = 5. Given the ground-truth set of target files \mathcal{F} , we define the reward as the proportion of correctly retrieved files:

Reward =
$$\frac{|\mathcal{Y} \cap \mathcal{F}|}{|\mathcal{F}|}$$

Given the prompt in the RL dataset, we let the model generate response and self-improve through this reward signal. We use the REINFORCE++ Hu et al. (2025) algorithm with a fixed learning rate of 1e-6 for the actor model. During training, we sample 8 rollouts per prompt. The training batch size is 64, and the rollout batch size is 256. The model is trained for 3 epochs, with a maximum prompt length of 16k tokens and generation length of 4k tokens. Additional hyperparameters include a KL divergence coefficient of 0.0, entropy coefficient of 0.001 and a sampling temperature of 1.0.

²https://github.com/sgl-project/sglang

D.4 RETRIEVAL REWARD MODEL

To train a reward model capable of reliably identifying the most relevant code files for modification, we construct a reward model dataset derived from our main dataset. The final reward model dataset consists of 112,378 samples corresponding to 25,363 unique instances. For each instance, the prompt is constructed using the retrieval reward model prompt template (see Appendix E), incorporating the issue statement along with the code content of each of the top-5 candidate retrieval files. Each data point is labeled with a binary value $\in 0, 1$, indicating whether the provided code content belongs to the ground-truth retrieval files. The model is initialized from the Qwen2.5-Coder-32B-Instruct and trained as a binary classifier using cross-entropy loss. Training is conducted with a batch size of 128, a learning rate of 5e-6, and a maximum sequence length of 32,768 tokens, over two epochs.

D.5 CODE EDITING MODEL

Dataset construction As described in Section 4, our code editing model is trained in two stages using supervised fine-tuning (SFT)—classical SFT and mutation SFT—followed by large-scale reinforcement learning (RL). We randomly select 1,470 instances (5%) from the full dataset for the classical SFT set and a separate 1,470 instances (5%) for the mutation SFT set. These two subsets are kept disjoint to ensure that the model learns to self-refine without direct exposure to the ground-truth solutions. For RL training, we use the remaining 22,102 instances (90% of the dataset), filtering out any prompts with sequence lengths exceeding 16,384 tokens. The RL dataset contains only the prompt (issue + code context) as input and the corresponding ground-truth patch as the output.

To synthesize the reasoning chain-of-thought (CoT) for classical SFT, we prompt Deepseek-V3-0324. Unlike the retrieval setting, we do not use rejection sampling, as Deepseek-V3-0324 often fails to generate the correct patch even after multiple samples. Instead, we adopt a more efficient approach by designing a "role-playing" prompt that provides the model access to the ground-truth patch and instructs it to explain the reasoning process behind it (see the "Generating Reasoning CoT for Code Editing Model (Classical SFT)" prompt in Appendix E). This ensures that the generated reasoning is both accurate and reflects an independent thought process. We then synthesize the classical SFT dataset using the "Code Editing Model (Classical SFT)" prompt template in Appendix E.

We first fine-tune the base model on the classical SFT dataset. This fine-tuned model is then used to generate five random patch candidates per instance with a sampling temperature of 1.0. These candidate patches are used to construct the mutation SFT dataset. For each instance, we prompt Deepseek-V3-0324 with: the issue statement, the content of the target file, the five candidate patches, and the ground-truth patch. Using the "Generating Reasoning CoT for Code Editing Model (Mutation SFT)" prompt (see Appendix E), the model is instructed to review each patch, critique their strengths and weaknesses, and propose an improved solution. We then extract the reasoning process and synthesize the mutation SFT dataset using the "Code Editing Model (Mutation SFT)" prompt template.

Supervised Fine-Tuning We perform supervised fine-tuning (SFT) on the Qwen2.5-Coder-32B-Instruct model using the synthetic SFT datasets described above. The prompt templates used for training are the same as those used to construct the two-stage SFT datasets (classical and mutation SFT). We employ a cosine learning rate scheduler with an initial learning rate of 1e-5. Training is conducted for one epoch, with a batch size of 128 and a maximum sequence length of 32,768 tokens.

Reinforcement Learning We fine-tune the mutation SFT model on the full dataset using REIN-FORCE++ Hu et al. (2025) with the following reward function:

$$r = \underbrace{R(x,y)}_{\text{Bonus}} + \underbrace{R(x,y) - \sum_{i=1}^{K} R(x,\bar{y}_i) - \lambda F(y)}_{\text{Format}}, \tag{6}$$

where each term is defined as follows:

• R(x,y) (Bonus): Encourages the model to produce high-reward outputs. Although similar in effect to the potential term, including this bonus stabilizes training and consistently improves the model's average reward.

- $R(x,y) \sum_{i=1}^{K} R(x,\bar{y}_i)$ (Potential): Measures the improvement of the current patch y over the average reward of the K conditioning patches \bar{y}_i . See Section 4.3 for details.
- F(y) (Format): Penalizes outputs that violate format or syntax constraints. It consists of:
 - String matching: Rewards outputs that closely match the ground-truth patch y^* using sequence similarity, following Wei et al. (2025).
 - Syntax check: Ensures the output can be parsed into the expected search-replace format, passes Python's ast syntax check, and satisfies flake8 static analysis. If any check fails, the format reward is set to zero.

The RL model is trained on a mix of data with and without conditioning examples. Conditioning examples are generated not only using the classical SFT model but also using a checkpoint of an RL-trained model at the first epoch.

As for implementation, we use REINFORCE++ Hu et al. (2025) algorithm with a fixed learning rate of 1e-6 for the actor model. During training, we sample 8 rollouts per prompt. The training batch size is 64, and the rollout batch size is 256. The model is trained for only 1 epochs, with a maximum prompt length of 16k tokens and generation length of 8k tokens. Additional hyperparameters include a KL divergence coefficient of 0.0, entropy coefficient of 0.001 and a sampling temperature of 1.0.

D.6 CODE EDITING REWARD MODEL

The code editing reward model is designed to provide a more accurate reward signal, addressing the limitations of using simple string-matching scores. The training setup is similar to that of the retrieval reward model (see Appendix D.4), with the main difference being in the data collection process. We construct the reward model training dataset using data collected from nebius/SWE-agent-trajectories and nebius/SWE-bench-extra, resulting in 56,797 samples across 1,889 unique instances. For each instance, the prompt is constructed using the code editing reward model prompt template (see Appendix E), and includes the issue statement, the code content of the target file to be modified, and a candidate patch. Each sample is labeled with a binary value \in 0, 1, indicating whether the candidate patch successfully resolves the issue. The model is trained as a binary classifier using the same training settings as the retrieval reward model.

D.7 REPRODUCTION TEST GENERATOR

Following a similar approach to that used for code editing, we generate intermediate reasoning steps for reproduction test generation using the Deepseek-V3-0324 model. Given the issue description and the corresponding ground-truth test patch, the model is prompted to produce a response that includes the reasoning behind constructing a valid test in a chain-of-thought format.

To support automated verification, we follow the strategy used in Agentless Xia et al. (2024), employing a test script template that prints clear diagnostic messages indicating whether the issue has been successfully reproduced or resolved. Specifically, If the test successfully triggers the target error (e.g., raises an AssertionError), it prints "Issue reproduced"; If the test completes without errors, it prints "Issue resolved". An example template for this diagnostic test script is shown below:

```
1890
         Reproduction Test Template
1891
1892
         def test_<meaningful_name>() -> None:
1893
              try:
1894
                  # Minimal code that triggers the bug
1895
              except AssertionError:
1896
                  print("Issue reproduced")
1897
                  return
1898
              except Exception:
1899
                  print("Other issues")
1900
                  return
1901
1902
              print("Issue resolved")
1903
              return
1904
1905
         if __name__ == "__main__":
              test_<meaningful_name>()
1906
1907
```

Starting from our filtered dataset, we generate one response per instance using greedy decoding and three additional responses via sampling with a temperature of 0.7. These synthetic examples are then used to fine-tune the <code>Qwen2.5-Coder-32B-Instruct</code> model over three epochs, resulting in our reproduction test generation model. The prompt templates used for generating intermediate reasoning and for supervised fine-tuning are provided in Appendix E.

E PROMPT TEMPLATE

```
1946
               Prompt Template — Retrieval Model
1947
1948
               Please look through the following GitHub problem description and Repository structure
               Determine the files most likely to be edited to fix the problem. Identify 5 most important files.
1949
               ### GitHub Problem Description ###
1950
               {problem_statement}
1951
               ### Repository Structure ###
1952
               {structure}
1953
               ### Format Instruction ###
1954

    Enclose reasoning process within `<think>...</think>`.

                   Please only provide the full path and return 5 most important files. Always return exactly 5 files,
1955
              Do Not output less than 5 or more than 5 files.

3. The returned files should be separated by new lines ordered by most to least important. Wrap all files together within `<file>...</file>`. 4. Do not include any explanations after `</think>`, only provide the file path within ` `<file>...</file>...
1956
1957
1958
               ### Examples ###
1959
               <think>
               1. Analyze the issue...
1960

    Check the files in provided repository structure for relevance...
    Confirm that the issue might be most relevant to 5 relevant files...

1961
               </think>
1962
               <file>
1963
               file1.py
               file2.pv
1964
               file3.py
1965
               file4.py
               file5.py
1966
               </file>
1967
1968
               Please provide your response below.
1969
1970
```

Prompt Template — Retrieval Reward Model

```
You are an expert software engineer and seasoned code reviewer, specializing in bug localization and 
→ code optimization.

You will be presented with a GitHub issue and a source code file.
Your task is to decide if the code file is relevant to the issue.

# Issue Statement
{problem_statement}

# File to be Modified
{file_content}
```

Prompt Template — Generating Reasoning CoT for Code Editing Model (Classical SFT)

```
You are a student striving to become an expert software engineer and seasoned code reviewer,

specializing in bug localization and code optimization within real-world code repositories. Your

strengths lie in understanding complex codebase structures and precisely identifying and modifying

the relevant parts of the code to resolve issues. You also excel at articulating your reasoning

process in a coherent, step-by-step manner that leads to efficient and correct

bug fixes.

You are now taking an exam to evaluate your capabilities. You will be provided with a codebase and an

issue description. Your task is to simulate a complete reasoning process—step-by-step—as if

solving the issue from scratch, followed by the code modifications to resolve the issue.

To evaluate your correctness, an oracle code modification patch will also be provided. You must ensure

that your final code modifications MATCH the oracle patch EXACTLY. However, your reasoning process

must appear fully self-derived and **must NOT reference, suggest awareness of, or appear to be

influenced by** the oracle patch. You must solve the problem as if you are unaware of the oracle

solution.

Issue Statement

{problem_statement}
```

```
1998
              Files to be Modified
1999
             Below are some code files that might be relevant to the issue above. One or more of these files may
             → contain bugs.
2001
             {file content}
2002
             # Oracle Code Modification Patch (For Evaluation Only):
2004
             {oracle patch}
2006
             # Reasoning Guidelines
2007
             Your reasoning process should generally follow these steps, with flexibility to adjust as needed for
            2008
             1. \star\starIssue Analysis\star\star: Start by thoroughly analyzing the issue. Explain what the problem is, why it
2009
                matters, and what the intended behavior should be. Identify the key goals and constraints that must
2010
            \hookrightarrow be addressed in your solution.
2011
            2. **Task Decomposition**: Break down the issue into smaller, manageable sub-tasks. Describe the
             purpose of each sub-task and how it contributes to solving the overall problem.
2012
             3. **Code Localization and Editing**: For each sub-task:
2013
                - Identify relevant code snippets by file path and code location.
2014
               - Explain how each snippet relates to the sub-task.
- Describe how the code should be changed and justify your reasoning.
2015
                - After thorough explanation, provide the corresponding edited code.
2016
             Your final output must precisely match the oracle patch, but your thinking must remain fully grounded
2017
            \hookrightarrow in the issue description and provided code files.
2018
2019
             # General Requirements
             1. **Independent and Evidence-Based Reasoning**: Your reasoning must be constructed as if independently
2020
            \hookrightarrow derived, based solely on the issue and code. Do not reference or imply knowledge of the oracle
2021

→ patch.

                **Clarity and Justification**: Ensure that each reasoning step is clear, well-justified, and easy to
2022
            \hookrightarrow
                follow.
             3. **Comprehensiveness with Focus**: Address all relevant components of the issue while remaining
2023
                concise and focused.
            4. **Faithful Final Output**: Your final code output must match the oracle patch exactly.
5. **Strict Neutrality**: Treat the oracle patch purely as a grading mechanism. Any hint of knowing the

→ patch in your reasoning (e.g., "based on the oracle," "we can verify," or "as we see in the patch")
2024
2025

→ will result in exam failure.

2026
2027
2028
            2029
2030
             3. Do not include any commentary or justification after the 

block.

2031
            Example:
2032
             <think>
             1. Analyze the issue...
2033
             2. Locate the relevant code...
             3. Apply necessary changes...
2034
2035
            ···python
2036
             # Final patch here (must match the oracle patch exactly)
2037
2038
2039
             Please provide your response.
2040
```

Prompt Template — Code Editing Model (Classical SFT)

20412042

2043

2044 2045

2046

2047

2048 2049

2050

```
You are an expert software engineer and seasoned code reviewer, specializing in bug localization and code optimization within real-world code repositories. Your strengths lie in understanding complex codebase structures and precisely identifying and modifying the relevant parts of the code to resolve issues. You also excel at articulating your reasoning process in a coherent, step-by-step manner that leads to efficient and correct bug fixes.

You will be provided with a codebase and an issue description. Your task is to simulate a complete reasoning process—step-by-step—as if solving the issue from scratch, followed by the code modifications to resolve the issue.

# Issue Statement
{problem_statement}
---
# Files to be Modified
```

```
2052
             Below are some code files that might be relevant to the issue above. One or more of these files may
2053
                contain bugs.
2054
             {file_content}
2055
             # Reasoning Guidelines
             Your reasoning process should generally follow these steps, with flexibility to adjust as needed for
2056
            1. **Issue Analysis**: Start by thoroughly analyzing the issue. Explain what the problem is, why it
2057
2058
            the intended behavior should be. Identify the key goals and constraints that must be addressed in your

→ solution.

             2. **Task Decomposition**: Break down the issue into smaller, manageable sub-tasks. Describe the
                purpose of each sub-task and how it contributes to solving the overall problem.
2060
             3. **Code Localization and Editing**: For each sub-task:
2061
                - Identify relevant code snippets by file path and code location.
                  Explain how each snippet relates to the sub-task.
2062
                  Describe how the code should be changed and justify your reasoning.
                - After thorough explanation, provide the corresponding edited code.
2063
2064
             # General Requirements
             1. **Clear and Evidence-Based Reasoning**: Provide clear and precise reasoning for each step, strictly
2065
            → based on the provided issue and code without inferring information not explicitly stated.
            2. **Comprehensive and Concise**: Address all relevant aspects of the issue comprehensively while being
2066

    ⇔ concise. Justify the exclusion of any sections that are not relevant.
    3. **Detailed Guidance**: Ensure the reasoning steps are detailed enough to allow someone unfamiliar

2067
                with the solution to infer and implement the necessary code modifications.
2068
             # Response Format
2069
             1. The reasoning process should be enclosed in <think> ... </think>.
            2. The final patch should be output in a standalone Python code block *after* the </think> block.
2070
             3. Do not include any commentary or justification after the </think> block.
2071
             # Patch Format
2072
             Please generate *SEARCH/REPLACE* edits to fix the issue. Every *SEARCH/REPLACE* edit must use this
                format:
2073
            1. The file path
             2. The start of search block: <<<< SEARCH
2074
                A contiguous chunk of lines to search for in the existing source code
2075
             4. The dividing line: ===
             5. The lines to replace into the source code
2076
             6. The end of the replace block: >>>>> REPLACE
             If, in 'Files to be Modified' part, there are multiple files or multiple locations in a single file
2077

→ require changes.

2078
            You should provide separate patches for each modification, clearly indicating the file name and the
                specific location of the modification.
2079
            Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION. For example, if you would like \hookrightarrow to add the line ' print(x)', you must fully write that out, with all those spaces before the \hookrightarrow code! And remember to wrap the *SEARCH/REPLACE* edit in blocks ```python...``
2080
             # Example Response
2081
             <think>
2082
             1. Analyze the issue..
             2. Locate the relevant code...
2083
             3. Apply necessary changes...
2084
               `pvthon
             ### mathweb/flask/app.py
2085
             <<<<< SEARCH
2086
             from flask import Flask
             import math
             from flask import Flask
             >>>>> REPLACE
2089
            ···python
2090
             ### mathweb/utils/calc.py
             <<<<< SEARCH
2091
             def calculate_area(radius):
                 return 3.14 * radius * radius
2092
2093
             def calculate area(radius):
                 return math.pi * radius ** 2
2094
             >>>>> REPLACE
2095
2096
            Please provide your response below.
2097
```

Prompt Template — Generating Reasoning CoT for Code Editing Model (Mutation SFT)

20982099

2100

2101

2102

2103

2104

2105

You are a student collaborating with a group of peers in a software engineering lab, working together \hookrightarrow to diagnose and fix bugs in real-world code repositories. You specialize in bug localization and \hookrightarrow code optimization, with a particular talent for critically evaluating others' patches and \hookrightarrow synthesizing high-quality, precise solutions from collaborative efforts.

You will be presented with a GitHub issue, the relevant source code files, and several *candidate \hookrightarrow patches* submitted by your teammates. Your task is twofold:

```
2106
             1. **Patch Review**: Carefully evaluate each of the several candidate patches **individually**.
2107
            \hookrightarrow Identify whether each patch resolves the issue correctly, partially, or incorrectly. If you identify any issues (e.g.,
2108
             → logical errors,
2109
            misunderstandings of the bug, overlooked edge cases, or incomplete fixes), explain them clearly and
             → suggest what could be improved or corrected.
2110
                Even if a patch appears mostly correct, you should still analyze its strengths and limitations in
2111
2112
                Treat this as a collaborative peer-review process: constructive, technical, and focused on improving

→ code quality.

2113
2114
            2. **Patch Synthesis**: After analyzing all several candidate patches, synthesize your understanding to

→ produce your **own final code patch** that fully resolves the issue. Your patch should:

2115
                  Be grounded solely in the issue description and provided source code.
                - Be informed by your peer review, but not copy any one patch outright.
2116
                 To evaluate your correctness, an oracle code modification patch will also be provided. You must
               \hookrightarrow ensure that your final code modifications MATCH the oracle patch EXACTLY. However, your
2117
               \hookrightarrow reasoning process must appear fully self-derived and **must NOT reference, suggest awareness of,
                   or appear to be influenced by** the oracle patch. You must solve the problem as if you are
2118
                  unaware of the oracle solution.
               \hookrightarrow
2119
2120
             # Issue Statement
2121
            {problem_statement}
2122
2123
            # Files to be Modified
2124
            Below are some code files that might be relevant to the issue above. One or more of these files may
2125

→ contain bugs.

2126
            {file content}
2127
2128
             # Candidate Patches (From Collaborators)
2129
            Below are several proposed patches submitted by your teammates. You will evaluate them individually.
             {candidate_patches}
2130
2131
            # Oracle Code Modification Patch (For Evaluation Only):
2132
            {target}
2133
2134
2135
            # Reasoning and Review Guidelines
2136
            Your response should be structured into two parts:
2137
             ## Part 1: Peer Patch Review
            For each of the candidate patches:
2138

    Analyze the candidate patch's intent and correctness.
    Identify what it does well, what it gets wrong (if anything), and how it could be improved.
    Use precise references to the provided issue and source code files to justify your evaluation.

2139
2140
                - Avoid any implication that you know the correct answer or are using an external reference
               2141
            ## Part 2: Final Patch Synthesis
2142
            After completing all reviews:
2143
            1. **Issue Analysis**: Start by thoroughly analyzing the issue. Explain what the problem is, why it
2144
                matters, and what the intended behavior should be. Identify the key goals and constraints that must

→ be addressed in your solution.

2145
            2. **Task Decomposition**: Break down the issue into smaller, manageable sub-tasks. Describe the
2146
            ← purpose of each sub-task and how it contributes to solving the overall problem.
2147
            3. **Code Localization and Editing**: For each sub-task:
2148
                - Identify relevant code snippets by file path and code location.
                 Explain how each snippet relates to the sub-task.
2149
                - Describe how the code should be changed and justify your reasoning.
2150
                - Incorporate useful insights from the candidate patches you reviewed. Reuse good ideas that are
                    correct and effective, but discard or correct those that contain flaws or misunderstandings.
2151
                - After thorough explanation, provide the corresponding edited code.
2152
            Your final output must precisely match the oracle patch, but your thinking must remain fully grounded

→ in the issue description and provided code files.

2153
2154
2155
            # General Requirements
            1. **Independent and Evidence-Based Reasoning**: Your reasoning must be constructed as if independently
2156

→ derived, based solely on the issue and code. Do not reference or imply knowledge of the oracle

2157
            2. **Clarity and Justification**: Ensure that each reasoning step is clear, well-justified, and easy to
2158

→ follow.
```

```
2160
             3. **Comprehensiveness with Focus**: Address all relevant components of the issue while remaining
2161

→ concise and focused.

2162
             4. **Faithful Final Output**: Your final code output must match the oracle patch exactly.
             5. **Strict Neutrality**: Treat the oracle patch purely as a grading mechanism. Any hint of knowing the \hookrightarrow patch in your reasoning (e.g., "based on the oracle," "we can verify," or "as we see in the patch")
2163
                 will result in exam failure.
2164
2165
2166
             # Response Format
             1. The reasoning process should be enclosed in <think> ... </think>.
2167
             2. The final oracle patch should be output in a standalone Python code block *after* the </think>
2168
             3. Do not include any commentary or justification after the </think> block.
2169
             Example:
2170
             <think>
             1. Review of candidate patch:
2171
                 - Review of patch-1: ..
                 - Review of patch-2: ...
2172
2173

    Analyze the issue by myself...
    Locate the relevant code...

2174
             4. Apply necessary changes...
             </think>
2175
2176
               `pvthon
             # Final patch here (must match the oracle patch exactly)
2177
2178
2179
             Please provide your response.
2180
```

Prompt Template — Code Editing Model (Mutation SFT)

2181

21822183

2184

2185

2186

2187

2188

2189

2190

2191

2192 2193

2194

219521962197

219821992200

2201

2202

2203 2204

2205

2206 2207 2208

2209

2210

2211

2212

```
You are an expert software engineer and seasoned code reviewer, specializing in bug localization and
    code optimization, with a particular talent for critically evaluating teammates' patches and
\hookrightarrow synthesizing high-quality, precise solutions from collaborative efforts.
You will be presented with a GitHub issue, the relevant source code files, and five *candidate patches*
\hookrightarrow submitted by your teammates. Your task is twofold:
1. **Patch Review**: Carefully evaluate each of the five candidate patches **individually**. Identify
\hookrightarrow whether each patch resolves the issue correctly, partially, or incorrectly. If you identify any \hookrightarrow issues (e.g., logical errors, misunderstandings of the bug, overlooked edge cases, or incomplete
\hookrightarrow fixes), explain them clearly and suggest what could be improved or corrected.
    Even if a patch appears mostly correct, you should still analyze its strengths and limitations in
   \hookrightarrow detail. Treat this as a collaborative peer-review process: constructive, technical, and focused \hookrightarrow on improving code quality.
2. **Patch Synthesis**: After analyzing all five candidate patches, synthesize your understanding to 

→ produce your **own final code patch** that fully resolves the issue. Your patch should:
      Be grounded solely in the issue description and provided source code.
    - Be informed by your peer review, but not copy any one patch outright.
# Issue Statement
{problem_statement}
# Files to be Modified
Below are some code files that might be relevant to the issue above. One or more of these files may

→ contain bugs.

{file_content}
# Candidate Patches (From Collaborators)
Below are five proposed patches submitted by your teammates. You will evaluate them individually.
{candidate_patches}
# Reasoning and Review Guidelines
Your response should be structured into two parts:
 ## Part 1: Peer Patch Review
For each of the five candidate patches:
      Analyze the candidate patch's intent and correctness.
    - Identify what it does well, what it gets wrong (if anything), and how it could be improved.

- Use precise references to the provided issue and source code files to justify your evaluation.
```

```
2215
             ## Part 2: Final Patch Synthesis
2216
             After completing all five reviews, your reasoning process should generally follow these steps, with
            \hookrightarrow flexibility to adjust as needed for clarity and accuracy:
2217
             1. **Issue Analysis**: Start by thoroughly analyzing the issue. Explain what the problem is, why it
2218
                 matters, and what the intended behavior should be. Identify the key goals and constraints that must
            \hookrightarrow be addressed in your solution.
2219
2220
            2. **Task Decomposition**: Break down the issue into smaller, manageable sub-tasks. Describe the
             → purpose of each sub-task and how it contributes to solving the overall problem.
2221
2222
             3. **Code Localization and Editing**: For each sub-task:
                  Identify relevant code snippets by file path and code location.
2223
                - Explain how each snippet relates to the sub-task.
                - Describe how the code should be changed and justify your reasoning.
2224
                - After thorough explanation, provide the corresponding edited code.
2225
2226
             # General Requirements
2227
             1. **Clear and Evidence-Based Reasoning**: Provide clear and precise reasoning for each step, strictly
                based on the provided issue and code without inferring information not explicitly stated
2228
            2. **Comprehensive and Concise**: Address all relevant aspects of the issue comprehensively while being 

→ concise. Justify the exclusion of any sections that are not relevant.
2229
            3. **Detailed Guidance**: Ensure the reasoning steps are detailed enough to allow someone unfamiliar
2230
            \hookrightarrow with the solution to infer and implement the necessary code modifications.
2231
2232
             # Response Format
2233
            1. The reasoning process should be enclosed in <think> \dots </think>.
             2. The final patch should be output in a standalone Python code block *after* the </think> block.
2234
             3. Do not include any commentary or justification after the </think> block.
2235
2236
             # Patch Format
2237
             Please generate *SEARCH/REPLACE* edits to fix the issue. Every *SEARCH/REPLACE* edit must use this

→ format:

2238
            1. The file path
             2. The start of search block: <<<< SEARCH
2239
             3. A contiguous chunk of lines to search for in the existing source code
2240
             4. The dividing line: ==
             5. The lines to replace into the source code
2241
             6. The end of the replace block: >>>>> REPLACE
2242
             If, in `Files to be Modified` part, there are multiple files or multiple locations in a single file
            \hookrightarrow require changes. You should provide separate patches for each modification, clearly indicating the \hookrightarrow file name and the specific location of the modification.
2243
2244
            Please note that the *SEARCH/REPLACE* edit REQUIRES PROPER INDENTATION. For example, if you would like \hookrightarrow to add the line ' print(x)', you must fully write that out, with all those spaces before the
2245
            ← code! And remember to wrap the *SEARCH/REPLACE* edit in blocks ```python...`
2246
             # Example Response
2247
             <think>
2248
             1. Review of candidate patch:
                - Review of patch-1: This patch attempts to fix X by modifying function Y. However, it fails to
2249
                - Review of patch-2: ...
2250
                  Review of patch-3: ...
2251
                - Review of patch-4: ...
                 Review of patch-5: ..
2252
             2. Analyze the issue by myself...
             3. Locate the relevant code..
2253
             4. Apply necessary changes...
             </think>
2254
2255
             ### mathweb/flask/app.py
2256
             <<<<< SEARCH
             from flask import Flask
2257
2258
             import math
             from flask import Flask
2259
             >>>>>> REPLACE
2260
2261
               `pvthon
             ### mathweb/utils/calc.py
2262
             <<<<< SEARCH
             def calculate area(radius):
2263
                 return 3.14 * radius * radius
2264
             def calculate area(radius):
2265
                 return math.pi * radius ** 2
             >>>>> REPLACE
2266
2267
```

```
Prompt Template — Code Editing Reward Model

* \begin{lstlisting} [language=text, fontsize=\tiny]
You are an expert software engineer and seasoned code reviewer, specializing in code optimization

→ within real-world
code repositories.
Your strengths lie in precisely identifying and modifying the relevant parts of the code to resolve
```

You will be provided with an issue description and an original code which has bugs. Your task is to write s code modifications to resolve the issue.

Problem Statement:
{problem_statement}

Original Code:

{file_content}minted % \end{lstlisting}

2283228422852286

226822692270

22712272

2273 2274

2275

2276

22772278

2279 2280

2281

2282

Prompt Template — Generating Reasoning CoT for Reproduction Test SFT

```
2287
              You are collaborating with peers in a software-engineering lab to create reproduction tests for

→ real-world bug reports.

2288
              You are given three context blocks:
2289
                -- BEGIN ISSUE (authoritative bug description) ---
2291
                - END ISSUE ---
                -- BEGIN ORIGINAL TEST FILES (do **not** reproduce the bug) ---
              {original tests}
2293
               -- END ORIGINAL TEST FILES ---
2294
                -- BEGIN TEST PATCH (contains a working reproduction) ---
2295
              {test_patch}
                -- END TEST PATCH ---
2296
2297
                **Important**
              > • The *Test patch* demonstrates at least one valid way to reproduce the bug; silently use it as
2298

→ inspiration to craft your own concise, single-file reproduction test.

> · **In your reasoning, act as if you derived everything from the Issue description alone.**

2299
              > • Do **not** refer to or hint at the presence of *Test patch*, *Original tests*, or any hidden
                   oracle.
              > • Your final script must follow the exact format below and reproduce *only* the behavior described in
2301
             \hookrightarrow the Issue.
2302
              ## Task
             Produce **one** self-contained Python test file that:
2305
              1. **Reproduces _only_ the bug described in the Issue** when the bug is present.
2. **Passes** (prints `"Issue resolved"`) once the bug has been fixed.
2306
              3. Prints exactly one of:
                 * "Issue reproduced" - bug still present (via AssertionError)

* "Issue resolved" - bug fixed / expectations met

* "Other issues" - unexpected exception unrelated to the Issue
2308
2309
              Reuse helpers from *Original tests* only if indispensable; otherwise keep the script standalone and
2310
2311
2312
              ## Response Format (**strict**)
2313
              1. Wrap **all reasoning** in a `<think> ... </think>` block.
2314
                 *Inside `<think>* you may explain how you interpreted the Issue and designed the test **without**
2315
                 \hookrightarrow mentioning or implying knowledge of the Test patch or any oracle.
2316
              2. After `</think>`, output **only** the final test script in a single Python code block.
2317
              Example skeleton * (follow this pattern exactly) *:
2318
2319
              <think>
              your independent reasoning here (no references to test_patch/oracle)
2320
              -
</think>
2321
```

```
2322
              `python
2323
            # All necessary imports
2324
            def test_<meaningful_name>() -> None:
2325
                   # minimal code that triggers the bug
2326
               except AssertionError:
2327
                   print("Issue reproduced")
2328
                   return
               except Exception:
2329
                   print("Other issues")
2330
                   return
2331
               print("Issue resolved")
               return
2332
2333
            if name == " main ":
                test_<meaningful_name>()
2334
2335
            **Guidelines**
2336
             **Focus solely on the Issue.** Strip out checks for any other problems that appear in *Test patch*.
2337
            * Keep the script **self-contained** unless a helper from *Original tests* is indispensable.
2338
            * Be concise--remove fixtures/parametrisations not strictly required.
2339
            Return your response in the exact format specified above.
2340
2341
```

```
2342
               Prompt Template — Reproduction Test Generator
2343
               You are collaborating with peers in a software-engineering lab to create reproduction tests for
2344
               \hookrightarrow real-world bug reports.
2345
               You are given the following authoritative bug description:
2346
                 -- BEGIN ISSUE --
2347
               {problem statement}
2348
                   - END ISSUE -
2349
                 **Important*

    You must independently derive a minimal reproduction test from the Issue description alone.
    Do **not** assume access to any "oracle," prior test patch, or original test files.
    Your final script must be self-contained and focused only on the behavior described in the Issue.

2350
2351
2352
2353
               ## Task
2354
               Produce **one** standalone Python test file that:
2355
               1. **Reproduces _only_ the bug described in the Issue** when the bug is present. 2. **Passes** (prints `"Issue resolved"`) once the bug has been fixed.
2356
               3. Prints exactly one of:
2357
                   * ""Issue reproduced" - bug still present (via AssertionError)

* ""Issue resolved" - bug fixed / expectations met

* "Other issues" - unexpected exception unrelated to the Issue
2358
2359
2360
               ## Response Format (**strict**)
2361
2362

    Wrap **all reasoning** in a `<think> ... </think>` block.

                   *Inside `<think>* explain how you interpreted the Issue and designed the test **without referencing
2363
                  \hookrightarrow any hidden tools, patches, or external files.***
2364
               2. After `</think>`, output **only** the final test script in a single Python code block.
2365
               Example skeleton \star (follow this pattern exactly) \star:
2366
               <think>
2367
               your independent reasoning here (no references to other tests or oracles)
               </think>
2368
2369
               # All necessary imports
2370
               def test_<meaningful_name>() -> None:
2371
                        # minimal code that triggers the bug
2372
2373
                    except AssertionError:
                        print("Issue reproduced")
2374
                         return
2375
                    except Exception:
```

```
2376
                    print("Other issues")
2377
                    return
2378
                print("Issue resolved")
2379
                return
2380
            if __name__ == "__main__":
    test_<meaningful_name>()
2381
2382
2383
            Guidelines
2384
2385
            Focus solely on the Issue description. Do not infer details not explicitly stated.
2386
            Keep the script self-contained--do not rely on external helpers or fixtures.
2387
            Be concise--remove all non-essential code and boilerplate.
2388
2389
2390
```